# Parallel Structure from Motion for UAV Images via Weighted Connected Dominating Set

San Jiang, Qingquan Li, Wanshou Jiang, Wu Chen

*Abstract*—Incremental Structure from Motion (ISfM) has been widely used for UAV image orientation. Its efficiency, however, decreases dramatically due to the sequential constraint. Although the divide-and-conquer strategy has been utilized for efficiency improvement, cluster merging becomes difficult or depends on seriously designed overlap structures. This paper proposes an algorithm to extract the global model for cluster merging and designs a parallel SfM solution to achieve efficient and accurate UAV image orientation. First, based on vocabulary tree retrieval, match pairs are selected to construct an undirected weighted match graph, whose edge weights are calculated by considering both the number and distribution of feature matches. Second, an algorithm, termed weighted connected dominating set (WCDS), is designed to achieve the simplification of the match graph and build the global model, which incorporates the edge weight in the graph node selection and enables the successful reconstruction of the global model. Third, the match graph is simultaneously divided into compact and non-overlapped clusters. After the parallel reconstruction, cluster merging is conducted with the aid of common 3D points between the global and cluster models. Finally, by using three UAV datasets that are captured by classical oblique and recent optimized views photogrammetry, the validation of the proposed solution is verified through comprehensive analysis and comparison. The experimental results demonstrate that the proposed parallel SfM can achieve 17.4 times efficiency improvement and comparative orientation accuracy. In absolute BA, the geo-referencing accuracy is approximately 2.0 and 3.0 times the GSD (Ground Sampling Distance) value in the horizontal and vertical directions, respectively. For parallel SfM, the proposed solution is a more reliable alternative.

*Index Terms*—unmanned aerial vehicle, structure from motion, 3D reconstruction, image orientation, bundle adjustment, connected dominating set

## I. INTRODUCTION

UAV (Unmanned aerial vehicle) has become one of the widely used remote sensing (RS) platforms in the field of photogrammetry and remote sensing. Compared with satellite and aerial RS platforms, UAVs have the characteristics of high flexibility, high timeliness, and high resolution [1], which have been used in transmission line inspection [2], [3], precision agriculture management [4], and cultural heritage document [5]. Efficient and accurate image orientation plays a critical role to ensure their applications in different domains.

Nowadays, SfM (Structure from Motion) has become the core technique for UAV image orientation [6]. Different from traditional POS (Positioning and Orientation System) aided aerial triangulation (AT), SfM can resume camera poses and scene 3D points from overlapped and unordered images without good initial values of unknown parameters [7], which has been implemented in well-known software packages as the AT module. According to the strategy for parameter initialization, SfM can be divided into three major groups, i.e., incremental SfM, global SfM, and hybrid SfM [8]. Compared with the other methods, incremental SfM has the advantages of robustness to outliers and high orientation precision, which is achieved through increasingly registering images for parameter initialization and iteratively executing bundle adjustment (BA) for parameter optimization [9]. This strategy, however, sacrifices the efficiency of image orientation. With the increase in data volumes and spatial resolutions, the capability for efficient and accurate orientation of UAV images becomes more and more important for modern SfM systems.

In the literature, different techniques have been designed to address the issues in the SfM-based image orientation pipeline. On the one hand, some research focused on the acceleration of BA optimization, including GPU (Graphics Processing Unit) and CPU (Central Processing Unit) based hardware acceleration [10] and the simplification [11] and optimization [12] of the BA mathematical model. However, these solutions cannot address the inherent drawback caused by the sequential constraint of incremental SfM. On the other hand, the divide-and-conquer strategy has been extensively adopted to break down the sequential constraint of incremental SfM [13]. The core idea is to divide the large-scale orientation problem into small-size and easy-to-process sub-problems, and the entire model is created by merging sub-models that can be reconstructed in an efficient distributed or parallel mode. Sub-model merging becomes the main difficulty. Simple methods utilize common camera poses or scene 3D points between sub-models, and the hierarchical methods organize sub-models as a tree structure and achieve model merging from bottom leaf nodes to the top root node. Their performance depends on the overlap structures between sub-models. For reliable model merging, global model constrained methods have also been documented, e.g., MLST (maximal leaf spanning tree) [14] and MCSD (minimal connected dominating set) [15]. These methods, however, are designed for landmark images with extremely high redundancy and overlap degrees.

S. Jiang and W. Chen are with the Department of Land Surveying and Geo-Informatics, The Hong Kong Polytechnic University, Hong Kong 999077, China; S. Jiang is also with the School of Computer Science, China University of Geosciences, Wuhan 430074, China. E-mail: *jiangsan@cug.edu.cn*, *wu.chen@polyu.edu.hk*.

Q. Li is with the College of Civil and Transportation Engineering, Shenzhen University, Shenzhen 518060, China, and also with the Guangdong Laboratory of Artificial Intelligence and Digital Economy (Shenzhen), Shenzhen 518060, China. E-mail: *liqq@szu.edu.cn*.

W. Jiang is with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan 430072, China. E-mail: *jws@whu.edu.cn*.

*Corresponding author: Wu Chen*

Considering the parallelism characteristic of the divide-and-conquer strategy, this study proposes a parallel incremental SfM solution for UAV images. The core idea of the proposed solution is to extract a global model from the image topological connection network (TCN), which is utilized as the global geometric constraint to assist model merging. Our main contributions are summarized as follows: (1) We propose a weighted connected dominating set (WCDS) algorithm to extract the global model, in which both node redundancy and edge weight are considered for TCN vertex selection. The proposed algorithm can enhance the edge connection of the global model and increase the completeness of SfM reconstruction. (2) We implement a workflow to achieve parallel SfM reconstruction, in which the image TCN graph can be established efficiently through the vocabulary tree-based image retrieval technique, and common 3D points between sub-models can be found efficiently through on-demand correspondence graph generation. (3) We verify the validation and demonstrate the performance of the proposed parallel SfM solution by using UAV datasets that are captured by both classical oblique and recent optimized views photogrammetry.

This paper is organized as follows. Section II gives a literature review that relates to SfM efficiency acceleration. Section III presents the workflow and detailed procedure of the proposed parallel SfM solution. By using real UAV datasets, a comprehensive analysis and comparison are presented in Section IV. Finally, Section V presents the conclusions of this study and improment plans for future study.

## II. RELATED WORK

This study focuses on the efficiency improvement of incremental SfM. In the literature, existing solutions can be categorized into two major groups, i.e., BA acceleration methods, and divide-and-conquer methods. Thus, the literature review is conducted from these two aspects as presented in the following subsections.

### A. BA acceleration methods

Incremental SfM depends on the iterative local and global BA to increase the precision of newly registered images and decrease the accumulated error of the final model. Iterative BA is a time-consumption step that dramatically degenerates the performance of image orientation. Therefore, BA acceleration is the most direct solution for SfM acceleration, including hardware acceleration, BA model simplification, and BA model optimization.

**Hardware acceleration**. Hardware techniques have tremendous development in recent years, which have been exploited to improve the efficiency of bundle adjustment [10], [16]–[19]. In the work of PBA (parallel bundle adjustment), [10] proposed using both multicore CPU and multicore GPU to solve the BA problem with high parallelism capability, in which the matrix-vector production operation is restructured to adapt to the hardware parallelizable mechanism. In the work of [19], the preconditioned conjugate gradient (PCG) and GPU-based parallel computing techniques were utilized to implement an efficient BA algorithm that was used for efficient

orientation for UAV images. In conclusion, orders of speedup ratios can be achieved via recent high-performance computing systems.

**BA model simplification**. In contrast to hardware acceleration, BA model simplification is another commonly used strategy to decrease the computational complexity of BA problems. Existing solutions are usually achieved by either decreasing the number of camera parameters [11], [20] or the number of 3D point structures [21], [22]. In the work of [11], a simplified BA model, termed RBA (reduced bundle adjustment), was proposed to process multi-camera oblique photogrammetric images, in which the pose of oblique cameras was simplified as the constant relative poses between oblique and vertical images and the absolute pose of vertical cameras. RBA reduces the total number of camera parameters in the BA optimization problem. Other solutions attempt to reduce the number of 3D points involved in the BA optimization since the parameters of 3D points are extremely larger than that of camera poses. This strategy is implemented by selecting the most useful tracks for image orientation [23], [24] or merely optimizing camera parameters in the BA problem, which is termed the structure-less SfM technique [21].

**BA model optimization**. For image orientation, the BA optimization is usually modeled as a joint minimization of reprojection errors between predicted and observed points, and it is solved as a nonlinear least-square problem. The sparse property of the normal equation in BA optimization can be utilized to decrease the number of involved parameters. On the one hand, the sparse property between cameras and 3D points is exploited by consecutive solving camera poses and 3D points, such as the SBA (sparse bundle adjustment) software package [12]. On the other hand, the sparse connection between cameras has been further used to handle the block data of SBA efficiently, which was released in the sSBA (sparse SBA), and it outperforms SBA by an order of magnitude. Except for these two packages, g2o (general graph optimization) is another well-known general BA optimization package that decreases the computational complexity of traditional BA solvers by exploiting the structure features of the BA problem [25].

### B. Divide-and-conquer methods

BA acceleration cannot break down the sequential constraint of incremental SfM, and they depend on high hardware requirements for increasingly large-scale image orientation. Thus, divide-and-conquer methods have been proposed, which mainly consist of simple methods, hierarchical methods, and global model constrained methods. For these methods, pairwise matches are first structured as an undirected weighted match graph, in which vertices indicate images, and edges weighted by matches are added for matched image pairs.

**Simple methods**. In this category, the initial match graph is first divided into sub-clusters with compact edge connections and small image size, and cluster merging is then achieved by using common camera poses or 3D points between sub-models [13], [26]–[29]. In the work of [13], the initial match graph is cut into small clusters through the normalized-cut (NC) algorithm [30], in which edge weights are assigned as the similarity

score of vocabulary tree based image retrieval. After the image orientation of each component, the entire scene is created by merging sub-models through their epipolar relationships. In the work of [29], two constraints, termed size constraint and completeness constraint, are used to divide the match graph into clusters with desired image number inside each cluster and enough image overlap between two clusters. Instead of the NC algorithm for scene clustering, [27] proposed using the matrix band reduction (MBR) algorithm because it can generate clusters with equal size and compact structure. Considering that the order of cluster merging affects the completeness and accuracy of the final model, [26] proposed the DagSfM algorithm that incorporates image graph and cluster graph into the steps of scene clustering and cluster merging, respectively. To improve the robustness of scene clustering and merging, [28] developed an automatic and dynamic strategy to split match graphs and merge sub-models, as well as evaluation metrics to search unreliable models.

**Hierarchical methods**. Instead of simply dividing the match graph into clusters at one time, hierarchical methods adopt a tree structure to organize scene clusters from bottom leaf nodes to the top root node, and cluster merging is executed by sewing up nodes from lower levels to higher levels. Hierarchical methods have been used as early as the arising of the SfM technique [31]. In the latter work of [32]–[34], the hierarchical cluster strategy was proposed for computing structure and motion, in which image similarity is calculated by using both correspondence number and their spatial distribution. The proposed SfM solution has been implemented in the commercial software package, termed SAMANTHA . Similarly, [35] adopted the divide-and-conquer manner to recursively divide an SfM problem into sub-maps. Compared with simple methods, hierarchical methods avoid the difficulties of selecting rational cluster size and determining the merging sequence.

**Global model constrained methods**. Hierarchical methods inherently lack the global view of the reconstruction problem since they only iteratively merge sub-models from bottom nodes to the top root node. In the literature, some other research proposes merging clusters from the middle tree level or using the global model to restrict the clustering merging. In the work of [36], base clusters were defined as the middle tree nodes with a pre-defined number of images, and the subsequent cluster merging started from these base clusters instead of the classical methods that start from tree leaf nodes. Inspired by the preemptive matching in [37], [38] first reconstructed a coarse and global model by using features with large scales, and the orientation of the remaining images was achieved in a parallel way through the direct resection using existing 3D-2D correspondences. In the work of [15], a graph optimization algorithm based on MCDS (minimal connected dominating set) was designed to reduce the number of vertices in the match graph and decrease the complexity of the subsequent SfM reconstruction. The proposed solution was verified by using landmark images with high redundancy from internet communities.
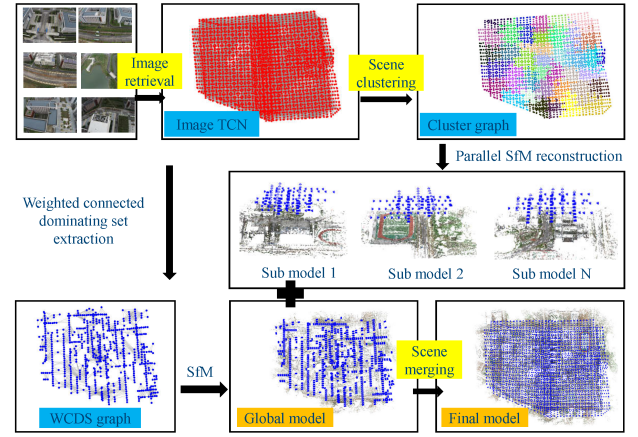


Fig. 1.  The workflow of the proposed parallel SfM solution.

## III. PARALLEL STRUCTURE FROM MOTION

Figure 1 presents the workflow of the proposed parallel Structure from Motion solution. The inputs of the parallel SfM solution are only UAV images and a pre-trained vocabulary tree without any other auxiliary data. The workflow consists of three major steps. First, overlapped match pairs are selected through vocabulary tree-based image retrieval, and an image TCN is then constructed as an undirected weighted graph, which establishes image connection and plays as the basic structure for global model extraction and scene clustering. Second, an algorithm, termed WCDS, is proposed to extract the global graph from the image TCN, which would be reconstructed and used as the global geometric constraint to guide scene merging. Meanwhile, the image TCN is simultaneously divided into small and compact clusters that can be reconstructed efficiently and accurately. Third, by using the WCDS graph and divided clusters, parallel SfM reconstruction is executed to generate the global model and cluster models, which are merged to generate the final model. This section first describes the used incremental SfM engine and then presents the workflow of the proposed parallel SfM solution.

### A. Incremental structure from motion

The workflow of the used incremental Structure from Motion engine consists of two modules, i.e., correspondence matching and incremental reconstruction [9]. The former aims to search for accurate and enough correspondences between overlapped match pairs, which can recover the relative geometry of two images; the latter is used to resume and refine camera poses and scene 3D points by using pair-wise geometry and BA optimization. In the field of photogrammetry and remote sensing, local feature-based image matching has become the widely used technique for correspondence matching because of their invariance to image rotation, scale difference, and even the changes of view-points and illuminations. In this study, the Root-SIFT with the L1-norm distance metric [39] has been adopted for feature extraction and matching.

After correspondence matching, incremental reconstruction is executed to recover camera poses and 3D points. Similar to

BA optimization in photogrammetry, pair-wise feature matches are first tied to generate tracks, which corresponds to a set of matched feature points that see the same 3D location [24]. Incremental SfM is then started by selecting two seed images that have a large enough intersection angle and a sufficient number of well-distributed matches, and the global model is constructed by resuming their relative poses and 3D points. For unregistered images, an iterative procedure is executed by searching for the next-best image that observes the largest number of resumed 3D points and registering the image into the reconstructed model. Meanwhile, both local and global BA optimization are alternately executed to refine the poses of newly added images and decrease accumulated errors [14].

For local and global BA optimization, the problem for refining camera poses and scene 3D points is formulated as a joint minimization of the reprojection function [40], where the sum of errors between the track projections and their corresponding image points is minimized. The object function of BA optimization is presented by Equation 1

$$\min_{C_j, X_i} \sum_{i=1}^{n} \sum_{j=1}^{m} \rho_{ij} \parallel P(C_j, X_i) - x_{ij} \parallel^2 \qquad (1)$$

where $X_i$ and $C_j$ indicate a 3D point and a camera, respectively; $P(C_j, X_i)$ is the projection of point $X_i$ on camera $C_j$; $x_{ij}$ is an observed image point; $\parallel \bullet \parallel$ denotes L2-norm; $\rho_{ij}$ is an indicator function with $\rho_{ij} = 1$ if point $X_i$ is visible in camera $C_j$; otherwise $\rho_{ij} = 0$. In this study, the Ceres Solver package is used for BA optimization.

### B. Match graph construction via vocabulary tree-based image retrieval

Match graph is used as the basic structure for clustering and merging in the parallel SfM solution, as well as guiding feature matching for efficiency improvement [41]. Before match graph construction, match pair selection should be conducted by using an efficient strategy. For most photogrammetric systems, onboard POS data is widely used to predict match pairs, in which image footprints on an average elevation plane are generated through the spatial intersection, and overlapped match pairs are determined through the intersection test between image footprints [42]. Although high efficiency can be achieved, these methods, however, depend on the precision of onboard POS data and the elevation of test sites, and they cannot adapt to data acquisition in optimized views photogrammetry.

In this study, vocabulary tree-based image retrieval has been utilized to perform match pair selection. In contrast to other methods, vocabulary tree-based image retrieval does not rely on any other auxiliary data but the images themselves. The core idea of vocabulary tree-based image retrieval is to represent each image as a BoW (Bag-of-Words) vector, and the problem of finding match pairs is cast as searching images with the closest BoW vectors between the query and database images [43], [44]. For match pair selection, a pre-trained and open-released vocabulary tree has been utilized in this study, which has 256 thousand visual words that are trained using images from internet communities. To achieve the highest efficiency, the number of used features for image indexing is set as 1500, and the number of retrieved images is configured as 100, which recovers enough match pairs.

After match pair selection, feature matching is then executed by using the SIFTGPU library for efficiency improvement, which is finally used for match graph construction. Suppose that $I = \{i_i\}$ and $P = \{p_{ij}\}$ are the images of size $n$ and match pairs of size $m$, respectively; the match graph is represented by an undirected graph $G = (V, E)$, where $V$ and $E$ stand for the vertex set and edge set of the undirected graph $G$, respectively. Thus, the match graph is constructed as follows: define a vertex $v_i$ for each image $i_i$, such that $V = \{v_i, i = 1, 2, ..., n\}$; an undirected edge $e_{ij}$ that connects vertices $v_i$ and $v_j$ is added for each of the $m$ match pairs $p_{ij}$, such that $E = \{e_{ij}, i, j = 1, 2, ..., n, i \neq j\}$. In addition, a weight value $w_{ij}$ is assigned to the edge $e_{ij}$.

In the context of SfM reconstruction, the weight value $w_{ij}$ encodes the connection strength of the corresponding edge $e_{ij}$. As feature matching has been conducted, it is rational to calculate the edge weight by using the number of feature matches. For image orientation, both the number and distribution of feature matches affect the robustness and precision of two view geometry estimation. Thus, the weight value $w_{ij}$ is calculated according to Equation 2

$$w_{ij} = R_{ew} * w_{inlier} + (1 - R_{ew}) * w_{overlap} \qquad (2)$$

where $w_{inlier}$ and $w_{overlap}$ are respectively weight items calculated by the number of feature matches and the overlap ratio between the convex of feature matches and the area of image planes; $R_{ew}$ is the weight ratio of these two items, which is set as 0.5 in this study. Thus, the weight value $w_{ij}$ is a linear combination of these two items, which are calculated according to Equations 3 and 4

$$w_{inlier} = \frac{log(N_{inlier})}{log(N_{max\_inlier})} \qquad (3)$$

$$w_{overlap} = \frac{CH_i + CH_j}{A_i + A_j} \qquad (4)$$

where $N_{inlier}$ and $N_{max\_inlier}$ respectively indicate the number of feature matches of the current match pair $p_{ij}$ and the maximum number of feature matches of all match pairs $P$; $CH_i$ and $CH_j$ are the convex hull areas of feature matches between images $i_i$ and $i_j$, respectively; $A_i$ and $A_j$ are the areas of image planes of images $i_i$ and $i_j$, respectively. In this study, the convex hull of matched feature points is detected by using the Graham-Andrew algorithm [45]. Noticeably, before the construction of the match graph, the match pairs with the number of matched features less than 50 are removed due to two main reasons. On the one hand, false matches may exist in these match pairs; on the other hand, they have a weaker connection in the subsequent global model construction.

### C. Parallel structure from motion

Parallel structure from motion can be implemented by dividing the original problem into some small-size and compact clusters that could be independently reconstructed based on the incremental SfM engine. In this study, both scene clustering

and cluster merging are achieved by using the match graph. First, a weighted connected dominating set is extracted from the initial match graph, which consists of a well connected subset of vertices and is used to reconstruct a global model to assist cluster merging. Second, compact clusters are generated by splitting the match graph into some independent parts without overlap vertices, which can be reconstructed efficiently under the memory limitation of computers. Finally, scene merging is conducted to generate the final model using the common 3D points between the global model and cluster models. The details of each step are described in the following subsections.

*1) Weighted connected dominating set for the global model:* Cluster merging is a non-trivial task in the SfM based on the divide-and-conquer strategy. In this study, the core idea of the proposed solution is to create a global model over the entire scene, and subsequent cluster merging can be constrained with the global model. To construct the global model, a tradeoff between two contradictory conditions should be made. On the one hand, the global model should be created by using as least number of images as possible as that the lowest computational costs can be consumed; on the other hand, the images involved in the global model should have strong and enough connections to ensure successful SfM reconstruction, which in turn would increases the number of images.

For the first condition, the purpose of extracting a minimal subset of images is the same as the minimal connected dominating set (MCDS) algorithm [46]. Suppose that the match graph is represented by $G = (V, E)$. MCDS aims to find a minimal subset of vertices $V_m$ from $V$ in graph $G$, such that the vertex set $V_m$ meets two conditions: first, the graph $G_m$ deduced from the vertex set $V_m$ is connected, which means that all vertices in $V_m$ are tied together; second, for any vertex $v_i$ in $V$, it belongs to the vertex set $V_m$; otherwise, it is the neighbor of at least one vertex in the vertex set $V_m$. Figure 2 illustrates the procedure of extracting MCDS from the match graph, which consists of four major steps as follows:

1) Initialization. In the beginning, all vertices in $V$ are labeled in while color, which indicates that all vertices are not visited. Meanwhile, set the vertex that has the largest number of white neighbors as the current vertex $v^*$, which is rendered in red color in Figure 2(a).
2) Scan and mark vertices. For the current vertex $v^*$, scan it (label it in black color). At the same time, mark the neighbors of $v^*$ (label them in gray color), as shown in Figure 2(b).
3) Find the next-current vertex. Among all marked vertices (vertices in gray color), set the vertex that has the largest number of white neighbors as the next-current vertex $v^*$, as shown in Figure 2(b) and Figure 2(c).
4) Iteratively execution of steps (2) and (3) until no while vertices exist. The vertices in black color consist of the extracted MCDS, as shown in Figure 2(d).

In the MCDS algorithm, the purpose of extracting a minimal subset of vertices from the match graph is implemented by selecting the next-current vertex $v^*$ that has the largest number of white neighbors. This strategy, however, does not consider the second condition of creating the global model. That is,
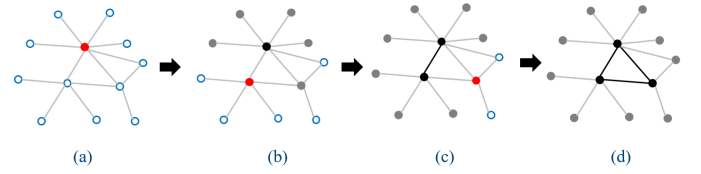


Fig. 2. The illustration of extracting MCDS: (a) initial match graph; (b) – (d) three vertices are iteratively appended into MCDS, and black nodes indicate the vertices of the extracted MCDS.

the images involved in the extracted vertex subset should have strong and enough connections since this condition ensures the precision and completeness of the final model in the context of SfM reconstruction. Thus, a weighted connected dominating set algorithm, termed WCDS, has been designed in this study. The core idea is to pose the edge weight constraint on the selection of the next-current vertex $v^*$. In particular, WCDS uses both the number of white neighbors and the edge weight between vertices to quantify the importance of gray vertices, as represented in Equation (5)

$$w_{ver} = R_{vw} * w_{ngb} + (1 - R_{vw}) * w_{ij} \quad (5)$$

where $w_{ngb}$ is the weight item calculated from the number of white neighbors; $w_{ij}$ is the edge weight that has been assigned to the corresponding edge $e_{ij}$ in match graph construction, as represented in Equation 2; $R_{vw}$ is the weight ratio that controls the contribution of these two weight items, which ranges between 0.0 and 1.0. Noticeably, when $R_{vw}$ is set as 1.0, the classical MCDS is utilized; when $R_{vw}$ is set as 0.0, only edge weight affects the calculation of the gray vertex importance. In Equation 5, the weight term $w_{ngb}$ has been normalized by Equation 6, in which $N_{ngb}$ is the number of white neighbors for the current vertex, and $N_{max\_ngb}$ is the maximal number of white neighbors for all vertices in the initial match graph as shown in Figure 2(a).

$$w_{ngb} = \frac{N_{ngb}}{N_{max\_ngb}} \quad (6)$$

Based on the weight calculation of gray vertices, the procedure of the proposed WCDS algorithm has been modified in step (3) when compared with the classical MCDS: for all marked vertices, calculate the gray vertex weight according to Equation 5, and set the vertex with the highest importance value $w_{ver}$ as the next-current vertex $v^*$. When one marked gray vertex connects to more than one scanned black vertex as shown in Figure 2(c), the vertex importance value of the marked gray vertex is calculated by using the largest edge weight. The main reason is that it is the strongest connection between the marked gray vertex and scanned black vertices.

*2) Scene clustering for parallel structure from motion:* Scene clustering aims to divide the entire scene into some small-size and compact clusters that can be efficiently and accurately reconstructed. To achieve the highest efficiency, each image that is not registered in the global model can be seen as an individual cluster, whose camera poses can be resumed efficiently using the 3D-2D correspondences deduced from

the reconstructed global model. However, considering that the global model only contains a very small fraction of images, it would be hard to search enough 3D-2D correspondences for unregistered images. Therefore, scene clustering is conducted to divide the match graph into some clusters, which can cover more 3D points in the global model.

Scene clustering has been implemented by splitting an initial match graph into some overlapped sub graphs under constraints, e.g., the size constraint and the completeness constraint. The former is used to control the size of each cluster for efficient parallel reconstruction; the latter is used to ensure enough common images between clusters for reliable model merging. Since the global model has been created to assist cluster merging, only the size constraint has been used in this study. For scene clustering, the normalized cut (NC) algorithm [30] has been selected, which is prone to cut graph edges with smaller weights and generate compact clusters with strong inner connections. After scene clustering, each image is assigned to only one cluster, and each cluster can be reconstructed based on the parallel SfM engine.

*3) Scene merging with the global geometric constraint:*
After the global model and cluster models are reconstructed based on the incremental SfM engine, each reconstructed model has its own coordinate system. The reconstructed models should be merged into the same coordinate system to generate a complete model. In general, scene merging is achieved through common structures between two cluster models, e.g., camera poses and 3D points. Suppose that two models $m_s$ and $m_r$ are respectively defined as the source reconstruction and reference reconstruction; there are $n$ common structures between these two models, as represented by $\{s_i\}$ and $\{r_i\}$, $i = 1, 2, ..., n$; the transformation from $m_s$ to $m_r$ is formulated by a similarity transformation $T = [\lambda R|t]$. In general, $T$ can be computed by minimizing the objective function as shown by Equation (7)

$$e^2(R, t, \lambda) = \frac{1}{n} \sum_{i=1}^{n} \parallel r_i - (\lambda R s_i + t) \parallel^2 \qquad (7)$$

where $R$ is the rotation matrix; $t$ is the translation vector; $\lambda$ is the scaling factor; $e^2$ is the mean square error (MSE) of transforming structures $\{s_i\}$ to $\{r_i\}$ under transformation $T$, which can be robustly estimated using the SVD (singular value decomposition) algorithm [47]. In this study, common 3D points between global and cluster reconstructions are used as the structure for similarity transformation estimation due to two main reasons. On the one hand, the number of 3D points is extremely larger than that of camera poses; on the other hand, 3D points are distributed more evenly in the global model.

Due to the existence of outliers in 3D points, the robust estimation technique based on the RANSAC framework has been utilized in similarity transformation estimation. Instead of using MSE for model verification, the bi-directional mean square reprojection error is utilized to calculate the transformation residuals under one hypothesis, as presented in Equation 8. In this study, the maximum residual $r_i$ is set as 1.8 pixels to separate inliers from outliers.
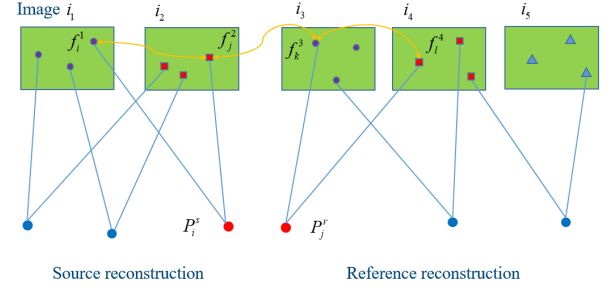


Fig. 3. The illustration of searching common 3D points by creating the on-demand correspondence graph between source and reference reconstructions.

$$r_i^2 = \frac{1}{m + l}(\sum_{j=1}^{m} \parallel P(C_j, T * s_i) - x_{ij} \parallel^2$$
$$+ \sum_{j=1}^{l} \parallel P(C_j, T^{-1} * r_i) - x_{ij} \parallel^2) \qquad (8)$$

Common 3D points between the global and cluster reconstructions should be efficiently determined during scene merging. In general, common 3D points can be found by either using common images between reconstructions [26] or by merging tracks from global and cluster reconstructions [29]. These methods either require high overlap degrees or need large memory consumption. In this study, an on-demand correspondence graph is created to establish the mapping relationship between feature matches, as illustrated in Figure 3. Two reconstructed models have respectively termed source and reference reconstructions. Within each reconstruction, the relationship between 2D feature points and 3D scene points has been established through the generated tracks in SfM reconstruction. In other words, the problem of finding common 3D points becomes establishing the mapping relationship between feature matches across these two reconstructions.

Considering that the number of images in the source reconstruction is much less than that in the reference reconstruction, the mapping relationship has been constructed by only using feature matches that are related to the images in the source reconstruction. As shown in Figure 3, only images labeled $i_1$, $i_2$, $i_3$ and $i_4$ are used to construct the correspondence graph, and the mapping relationship between feature matches are indicated by the yellow lines. Using the established correspondence graph, common 3D points $C$ are found according to the procedure: (1) for each 3D point $P_i^s$ in the source reconstruction, find its related observations $O = \{f_i^1, f_j^2\}$; (2) for each observation in $O$, by using the mapping relationship in the correspondence graph, find its related feature matches $M = \{f_k^3\}$ from the images in the reference reconstruction; (3) for each feature match in $M$, find its corresponding 3D points $P_j^r$ based on established tracks in the reference reconstruction; (4) add the common 3D point pair $(P_i^s, P_j^r)$ to $C$ if it does not exist in $C$. For scene merging, the proposed method based on on-demand correspondence graph has very high efficiency and low memory cost.

Two reconstructions can be merged by using their established similarity transformation within the robust RANSAC framework. Although all cluster reconstructions are transformed into the coordinate system of the base reconstruction, the merging sequence among the cluster reconstruction would affect the precision and completeness of the final model. Intuitively, the MSE from similarity transformation estimation may be a good metric to quantify the merging quality of two reconstructions, which has been used in other work [26]. This study, however, orders the merging sequence based on the number of common 3D points due to two main reasons. On the one hand, two reconstructions with more common 3D points are usually prone to be merged with high precision; on the other hand, the earlier more 3D points are merged into the global model, the more common 3D points can be found for the subsequent cluster models. After merging all cluster models, one final global BA is executed to optimize the camera poses and 3D points to refine the final model.

### D. Algorithm implementation

Based on the principle of the proposed solution, this study has implemented the parallel SfM system for UAV images. We have used the C++ programming language to achieve the highest efficiency. For feature detection and matching, the used SIFTGPU has been accelerated by using CUDA (Compute Unified Device Architecture), and the default parameters are used as documented in [48]. For SfM reconstruction, the open-source software package ColMap [49] has been used as the backend SfM engine. In this study, we have added the BROWN camera model that has been widely used in the compared commercial software packages and implemented an absolute BA optimizer to support geo-referencing accuracy assessment. In addition, camera intrinsic parameters have been calibrated and are fixed in BA optimization. The algorithmic implementation is presented in Algorithm 1

### IV. EXPERIMENTAL RESULTS AND DISCUSSIONS

In experiments, three UAV datasets are collected to evaluate the performance of the proposed parallel SfM solution. First, the influence of the weight ratio $Rvw$ on extracting the WCDS graph is analyzed in terms of the selected image ratio and connected image ratio. Second, memory consumptions are compared in the cluster merging with and without the on-demand correspondence graph strategy. Third, by using the optimal parameter setting, the SfM solution is executed to reconstruct the three UAV datasets. Finally, the proposed SfM solution is extensively compared with four open-source and commercial software packages, including ColMap [49], DagSfM [26], Agisoft Metashape [50], and Pix4dMapper [51], and the results are analyzed in relative BA without ground control points (GCPs) and absolute BA with GCPs.

### A. Test sites and datasets

Table I presents the detailed information of these three UAV datasets. For outdoor data acquisition, multi-rotor UAVs are used in the three campaigns. The details of each dataset are described as follows:

---

**Algorithm 1** Parallel SfM based on WCDS

**Input:** UAV images $I = \{i_i\}$; a pre-trained vocabulary tree $VocT$

**Output:** reconstructed model $M$

1: **procedure** TCNCONSTRUCTION
2:    **for each** $i_i \in I$ **do**
3:       Extract SIFT features for image $i_i$
4:       Select the top-scale 1500 features of image $i_i$ and index them into $VocT$
5:    **end for**
6:    **for each** $i_i \in I$ **do**
7:       Query the top 100 similar images and add them to match pairs $P = \{p_{ij}\}$
8:    **end for**
9:    Create $TCN$ graph using $P$ with edge weights calculated by Equation 2
10: **end procedure**

1: **procedure** WCDSEXTRACTION
2:    Initialize all vertices in $TCN$ as white color
3:    Select the vertex with the most white neighbors as current vertex $v^*$
4:    **do**
5:       Scan vertex $v^*$ and mark its neighbors
6:       Update the weight of all gray vertices according to Equation 5
7:       Set the gray vertex with the largest weight as next-current vertex $v^*$
8:    **while** white vertex exists
9:    Create the $WCDS$ graph using vertices in black color
10: **end procedure**

1: **procedure** PARALLELRECONSTRUCTION
2:    Start a thread to reconstruct $WCDS$ based on incremental SfM
3:    Divide $TCN$ into non-overlapped clusters $C = \{c_j\}$ of size $m$
4:    **for each** $c_j \in C$ **do**
5:       Start a thread to reconstruct $c_j$ based on incremental SfM
6:    **end for**
7: **end procedure**

1: **procedure** CLUSTERMERGING
2:    Set $M = WCDS$
3:    **while** $C$ not empty **do**
4:       Set $P = \{\}$
5:       **for each** $c_j \in C$ **do**
6:          Find common 3D points $p_j$ between $c_j$ and $M$, set $P = P + \{p_j\}$
7:       **end for**
8:       Find $p_k^*$ that has the largest number of common 3D points
9:       Calculate transform $T_k$ using RANSAC according to Equation 8
10:      Merge $c_k$ into $M = M + \{c_j\}$ and set $C = C - \{c_j\}$
11:    **end while**
12:    Conduct the final BA for the merge model $M$
13: **end procedure**
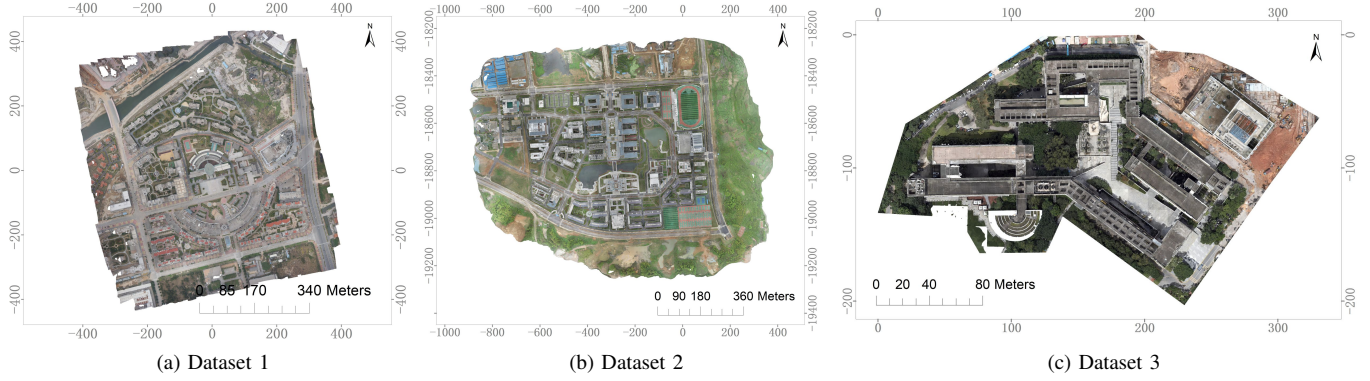
(a) Dataset 1      (b) Dataset 2      (c) Dataset 3

Fig. 4. The orthomosaics for the test sites of the three datasets.

- The first dataset covers an urban shopping plaza that is surrounded by high buildings, as shown in 4a. For image acquisition, a classical five-camera oblique photogrammetric system consisting of one nadir camera and four oblique cameras is utilized, where four oblique cameras are rotated 45° with respect to the nadir camera. The resolution of captured images is 6000 by 4000 pixels. Under the flight height of 175 m, a total number of 750 images are recorded with the GSD value of 4.3 cm.
- The second dataset covers a university campus, which includes dense low buildings, as shown in 4b. The surroundings of the campus are vegetation and bare land. By using a DJI Phantom 4 RTK UAV equipped with one DJI FC6310R camera, a total number of 3743 images with the dimension of 5472 by 3648 pixels are collected under the flight height of 80 m. The GSD of collected images is approximately 2.6 cm.
- The third dataset is recorded from a university campus, which is mainly covered by a complex building, as shown in 4c. In contrast to the other datasets, the optimized views photogrammetric technique [52] has been used for designing the UAV trajectory, which adjusts view points and directions according to the coarse geometric model of the test site. By using the DJI M300 RTK UAV equipped with one DJI Zenmuse P1 camera, a total number of 4030 images are recorded with the GSD of 1.2 cm under the view distance of approximately 80 m to ground targets. Besides, for the geo-reference accuracy assessment, 26

GCPs are surveyed using a total station, whose nominal horizontal and vertical accuracy is about 0.8 cm and 1.5 cm. GCPs are made by using artificial markers and distributed on ground, building facades and roofs.

### B. The influence of the weight ratio on the construction of the WCDS graph

In the construction of the WCDS graph, the weight ratio $R_{vw}$ controls the contribution of the number of white neighbors and the edge weight to quantify the importance of gray vertices. This section would analyze the influence of the weight ratio $R_{vw}$ on creating the WCDS graph.



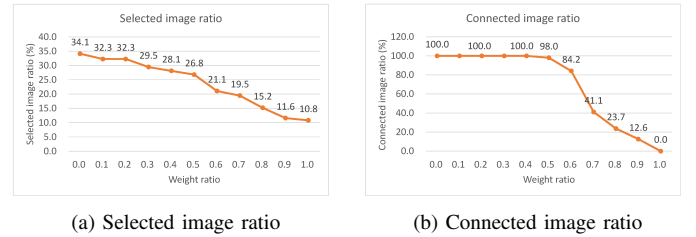(a) Selected image ratio      (b) Connected image ratio

Fig. 5. The influence of the weight ratio on global geometric constraint construction for dataset 1.

For the analysis of the weight ratio, dataset 1 is selected for experiments, and two metrics are utilized for performance evaluation, i.e., the selected image ratio in WCDS graph and the connected image ratio in SfM reconstruction. The former is the ratio of the number of selected images in the WCDS graph and the number of all images in the dataset; the latter is the ratio of the number of connected images in SfM reconstruction and the number of selected images in the WCDS graph. In this test, the weight ratio is uniformly sampled between 0.0 and 1.0 with the interval value of 0.1. The results are shown in Figure 5, where the influence on the selected image ratio is presented in Figure 5a and the influence on the connected image ratio is shown in Figure 5b. It is shown that with the increase of the weight ratio from 0.0 to 1.0, the selected image ratio is almost linearly decreasing from 34.1% to 10.8%. The main reason is that by using a smaller weight ratio, WCDS tends to select the next-current vertex that has a stronger edge connection, which further slowdowns the speed of vertex scanning and marking.

TABLE I
DETAILED INFORMATION OF THE THREE UAV DATASETS.

| Item Name | Dataset 1 | Dataset 2 | Dataset 3 |
|---|---|---|---|
| UAV type | multi-rotor | multi-rotor | multi-rotor |
| Flight height (m) | 175 | 80 | - |
| Camera mode | Sony NEX-7 | DJI FC6310R | DJI ZenmuseP1 |
| Camera number | 5 | 1 | 1 |
| Focal length (mm) | nadir: 16 oblique: 35 | 24 | 35 |
| Camera angle (°) | nadir: 0 oblique: 45/−45 | 35 | - |
| Number of images | 750 | 3743 | 4030 |
| Image size (pixel) | 6000 × 4000 | 5472 × 3648 | 8192 × 5460 |
| GSD (cm) | 4.3 | 2.6 | 1.2 |

(a) Original TCN graph     (b) WCDS graph with value 0.0     (c) WCDS graph with value 0.2

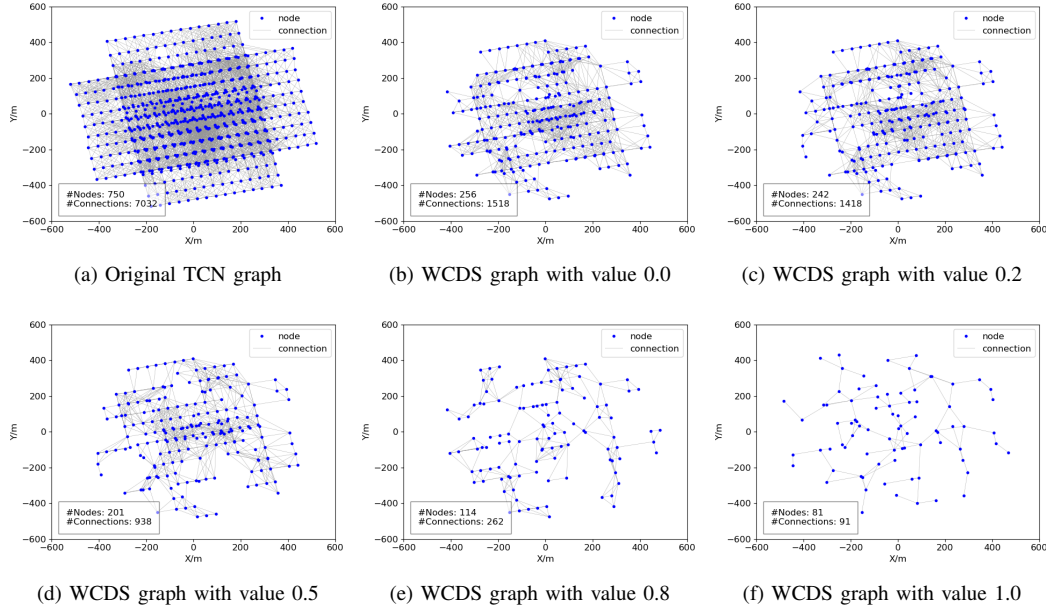(d) WCDS graph with value 0.5     (e) WCDS graph with value 0.8     (f) WCDS graph with value 1.0

Fig. 6. The illustration of the weight ratio on WCDS extraction for dataset 1.

By observing the results of the connected image ratio presented in Figure 5b, we can see that with the increase of the weight ratio from 0.0 to 0.5, almost all selected images have been successfully connected in SfM reconstruction. However, the connected image ratio decreases dramatically when the weight ratio decreases from 0.5 to 1.0. Noticeably, no images have been connected when the weight ratio is set as 1.0. In other words, the extracted images from MCDS has a very weak connection. For visual analysis, Figure 6 illustrates the weight ratio on the construction of the WCDS graph, in which Figure 6a shows the original TCN graph, and Figure 6b - Figure 6f show the WCDS graph with the weight ratio of 0.0, 0.2, 0.5, 0.8 and 1.0, respectively. It is clearly shown that the weight ratio has good control on the connection stability of the WCDS graph. By combining the results of the selected image ratio and the connected image ratio, we set the weight ratio $R_{vw}$ as 0.5 in the following tests.

## C. The influence of correspondence graph generation on cluster merging

In cluster merging, the correspondence graph is created to build the mapping relationship between feature matches and would be used to find common 3D points across reconstructions. With increasing the number of images, it becomes impossible to build tracks at once for the entire dataset due to the memory limitation and computational costs. In the proposed solution, the on-demand correspondence graph is built by using the feature matches between the images in the source reconstruction and the related images in the reference reconstruction. In this test, we would analyze the influence correspondence graph on cluster merging.

In this section, three solutions for building the correspondence graph have been analyzed, including all-cluster, pairwise-cluster, and on-demand strategies. The first one builds
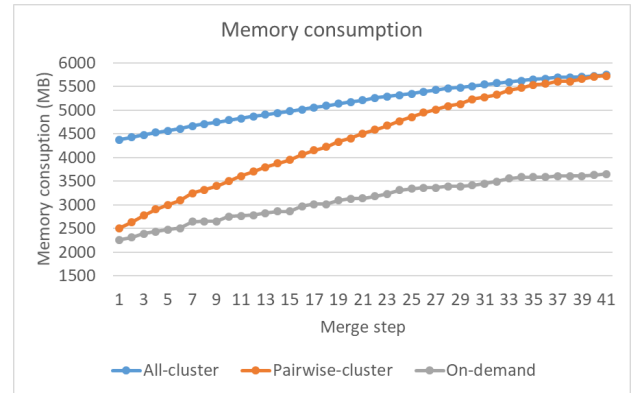


Fig. 7. The statistic result of memory consumption during cluster merging for dataset 3.

the correspondence graph at once using all feature matches from the dataset; the second one uses all feature matches from two reconstructions; the third one uses feature matches from related images between two reconstructions. For performance evaluation, memory consumption has been used as the metric. Figure 7 shows the statistic result in the cluster merging of dataset 3. It is shown that the all-cluster strategy needs the highest memory consumption as it builds the correspondence graph for the entire dataset; although the time consumption of the pair-wise strategy is the same as the proposed on-demand strategy at the beginning, it increases with the number of merged clusters and reaches to the same value as the all-cluster strategy. On the contrary, stable memory consumption can be observed from the on-demand strategy as it only loads required feature matches between reconstructions. For dataset 3, there is approximately 2.1 G memory saving during cluster merging. With the increase of the dataset, more memory saving can be achieved by using the on-demand strategy.
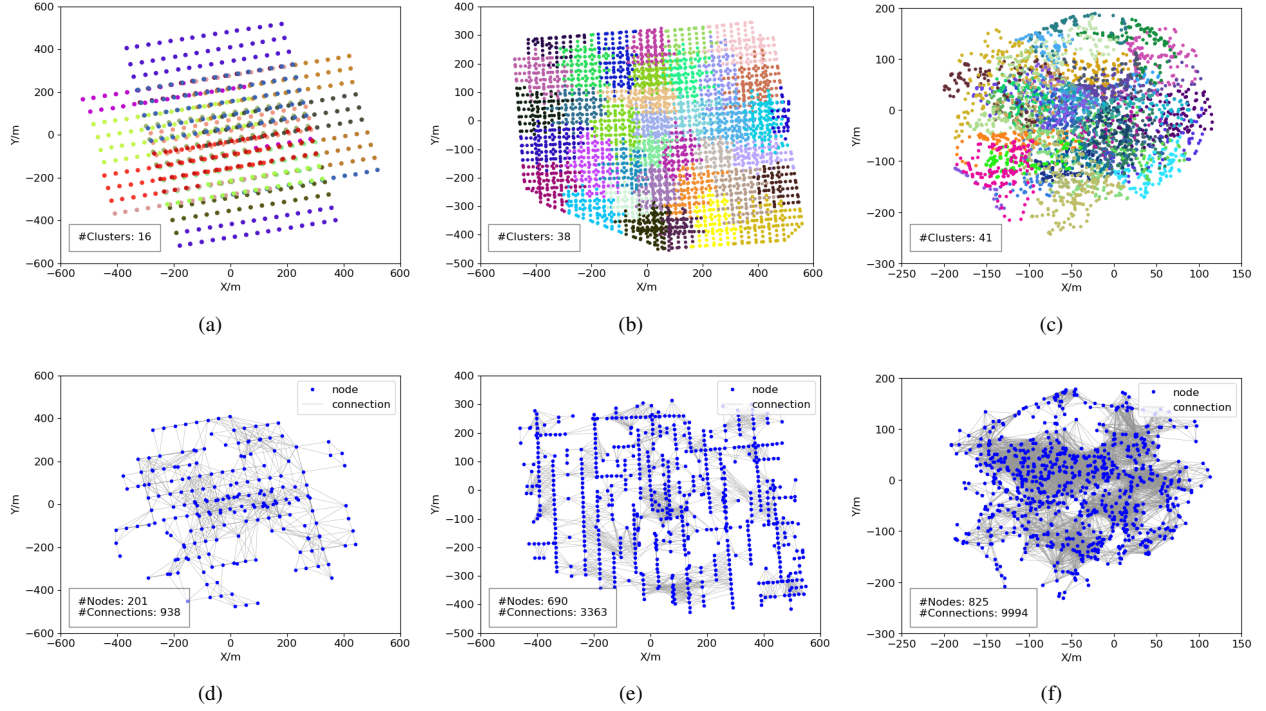
Fig. 8. The illustration of scene partition and weighted CDS for datasets 2 and 4: (a), (c), and (e) scene partition with the cluster size of 50, 100, and 100 for these three datasets, respectively; (b), (d) and (f) weighted CDS for the construction of global geometric constraint for the three datasets, respectively.

*D. Parallel structure from motion for test datasets*

By using the selected optimal parameters, the WCDS graph is extracted to build the global model, are shown in Figure 8d, Figure 8e, and Figure 8f. The number of selected images is 201, 690, and 825 for the three datasets, respectively, which is about 26.8%, 18.4%, and 20.5% of the total number of images in the corresponding dataset. We can also observe that enough image connections have been established in the extracted WCDS graph, which can be verified by the dense gray edges in each WCDS graph. In addition, Figure 9a, Figure 9c and Figure 9e present the SfM reconstructions based on the WCDS graphs, and almost all images in the corresponding graph have been successfully registered, whose number is 197, 684 and 821, respectively.

Meanwhile, scene cluster is also executed to divide the whole scene into clusters. For these three datasets, the number of cluster size is set as 50, 100, and 100, and the number of generated clusters is 16, 38, and 41, which are illustrated in Figure 8a, Figure 8b and Figure 8c, respectively. Noticeably, the cluster result shown in Figure 8c seems much more irregular than the results presented in Figure 8a, and Figure 8b. The main reason is that dataset 3 has been collected by using the optimized views photogrammetry, in which images are not captured at a fixed altitude. After the parallel SfM reconstruction guided by the clusters and the cluster merging constrained by the global model, the final models have been generated and presented in Figure 9b, Figure 9d, and Figure 9f for the three datasets, respectively. We can see that almost all images have been successfully connected for datasets 1 and 2 that are captured by classical oblique photogrammetry.

On the contrary, there are 123 images lost in the merged reconstruction of dataset 3, which may be caused by the lacking of enough common 3D points since the irregular data acquisition campaign has been utilized in this dataset.

*E. Comparison with state-of-the-art methods*

To assess the proposed parallel SfM solution, comparison tests with both open-source and commercial software packages have been conducted in terms of efficiency, precision, and completeness. First, relative BA without GCPs is performed to assess the reconstruction efficiency, relative precision, and completeness of the final model. Second, absolute BA with GCPs is conducted by using ground-truth data to assess the geo-reference accuracy of the final model. For the comparison, two open-source packages, termed ColMap and DagSfM, and two commercial packages, termed Agisoft Metashape and Pix4Dmapper, are used, and these two open-source software packages are implemented by using the C++ programming language. In addition, default parameters are configured for these packages, and all comparison tests are conducted on an Intel Core i7-8700 PC on the Windows platform with 32 GB memory, a 3.19 GHz CPU, and a 6 GB NVIDIAN GeForce GTX 1060 graphics card.

*1) Relative BA in terms of efficiency, precision, and completeness:* For relative BA without GCPs, three metrics have been used for performance evaluation, including efficiency, precision, and completeness. The metric efficiency indicates the time costs used in the parallel SfM reconstruction; the metric precision is calculated as the mean reprojection error after the execution of the global BA; the metric completeness
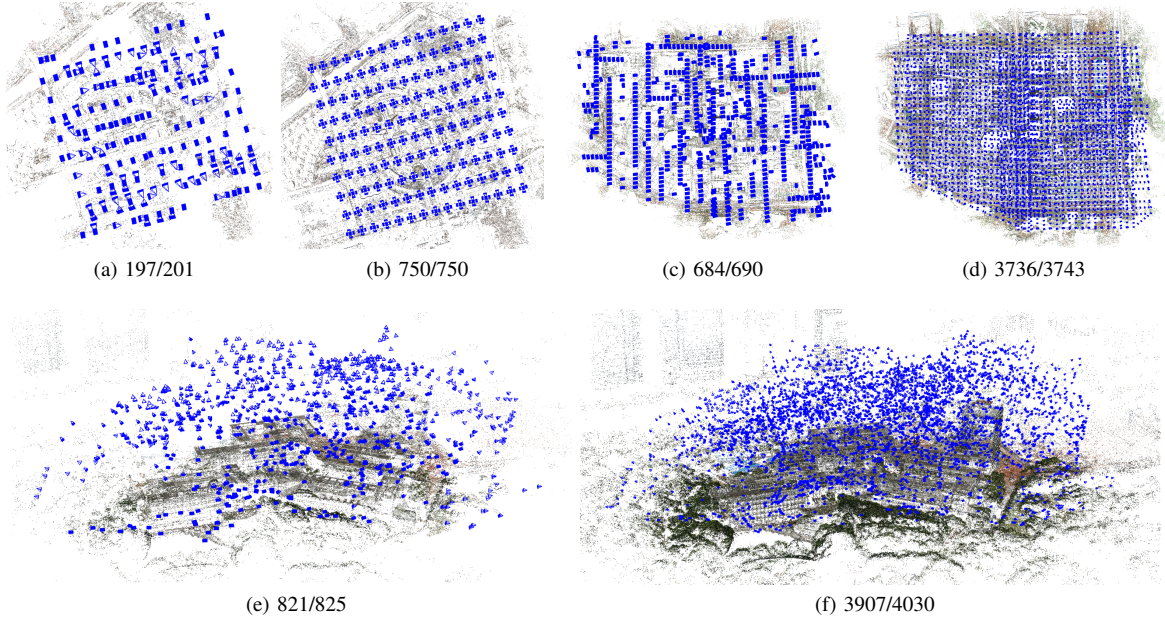
Fig. 9. The results of SfM-based 3D reconstruction of the three datasets: (a), (c), and (e) 3D reconstruction results of global geometric models; (b), (d), and (f) 3D reconstruction results of all images in the corresponding datasets. The values in each sub-figure indicate the numbers of resumed and contained images, respectively.

is quantified by using the number of registered images and resumed 3D points.

Table II presents the statistical results of relative BA for the three datasets. We can see that the proposed solution achieves the highest efficiency, which is 4.81 min, 100.63 min, and 186.96 min for the three datasets, respectively. The efficiency of DagSfM ranks second since it also utilizes a parallel SfM reconstruction method. When the number of images is not too large, the time costs for the other three packages are acceptable, such as for dataset 1, the time costs are 20.49 min, 43.00 min, and 10.88 min for ColMap, Metashape, and Pix4Dmapper, respectively. However, their efficiency degenerates dramatically with the increase of the involved images. Especially for ColMap, the time cost increases to 3,256.68 min for dataset 3 with 4030 images, which is approximately 17.4 times the time costs consumed by the proposed solution. Even worse, Metashape fails to reconstruct dataset 3 due to the reason being out of memory.

When considering the metrics precision, we can observe that Pix4Dmapper achieves the best performance among all compared packages, which is followed by the proposed solution with the value of 0.410 pixels, 0.374 pixels, and 0.429 pixels for the three datasets. Compared with the parallel method in DagSfM, the proposed solution utilizes the WCDS graph to provide the global geometric constraint in cluster merging and achieves higher BA precision. When considering the metric completeness, it is shown that the largest number of 3D points have been reconstructed from Metashape because it uses the divide-and-conquer strategy for feature detection. For dataset 1 with the smallest image number, all images can be registered by the compared packages. With the increase in the number of involved images, DagSfM lost many more images in the final model when compared with the other packages,

which can also be verified by the reconstructed point clouds of dataset 2 as presented in Figure 10. The main reason is that DagSfM merely depends on the overlapped images between reconstructions for cluster merging. By using the global models, the proposed solution registers 3,736 and 3,907 images for datasets 2 and 3, respectively. Based on the results of relative BA, we can conclude that by using the WCDS-based global geometric constraint, the proposed parallel SfM solution can achieve an order of efficiency improvement when compared with classical incremental SfM solution.

*2) Absolute BA in the term of geo-referencing accuracy:*
With ground-truth data, absolute BA is conducted to evaluate geo-referencing accuracy. During data acquisition of dataset 3, a total number of 26 GCPs have been surveyed and prepared for geo-referencing accuracy assessment. In this test, three GCPs that distribute evenly in the test site are selected to involve the absolute BA for geo-referencing, and all the others are used as check points (CPs). For the accuracy assessment, the residuals of model geo-referencing are calculated as the coordinate differences between CPs and their model points that are triangulated from the geo-referenced model.

Table III shows the statistical results of geo-referencing accuracy, and three metrics, termed max, mean and std.dev. of absolute residuals, are used for performance comparison. We can see that Pix4Dmapper achieves the highest accuracy among all compared methods, whose residual in the term of std.dev. is 0.013 m, 0.016 m, and 0.019 m in the X, Y, and Z directions, respectively. Although the vertical accuracy of ColMap ranks second with the std.dev. value of 0.018 m, the proposed solution can achieve comparative accuracy in the horizontal directions with the std.dev. value of 0.020 m, and 0.024 m in the X and Y directions, respectively, which can also be demonstrated by the individual residual plots as presented

TABLE II
THE STATISTICAL RESULTS OF RELATIVE BA WITHOUT GCPS FOR THE THREE DATASETS IN TERMS OF EFFICIENCY, COMPLETENESS, AND PRECISION.
THE VALUES IN THE BRACKET INDICATE THE NUMBER OF CONNECTED IMAGES IN THE FINAL MODELS.

| Metric | Method | Dataset 1 | Dataset2 | Dataset 3 |
|---|---|---|---|---|
| Efficiency (min) | ColMap | 20.49 | 387.84 | 3,256.68 |
| | DagSfM | 7.79 | 129.06 | 197.04 |
| | Metashape | 43.00 | 208.00 | — |
| | Pix4Dmapper | 10.88 | 290.72 | 753.90 |
| | Ours | 4.81 | 100.63 | 186.96 |
| Precision (pixel) | ColMap | 0.683 | 0.580 | 0.712 |
| | DagSfM | 1.000 | 0.722 | 0.739 |
| | Metashape | 0.562 | 0.889 | — |
| | Pix4Dmapper | 0.253 | 0.379 | 0.373 |
| | Ours | 0.410 | 0.374 | 0.429 |
| Completeness | ColMap | 308,670 (750) | 1,211,943 (3,738) | 1,594,456 (4,029) |
| | DagSfM | 252,919 (750) | 1,258,696 (3,254) | 1,510,435 (3,487) |
| | Metashape | 905,815 (750) | 7,733,683 (3,742) | — |
| | Pix4Dmapper | 627,919 (750) | 4,835,383 (3,743) | 9,987,579 (3,973) |
| | Ours | 231,780 (750) | 1,253,274 (3,736) | 1,557,104 (3,907) |



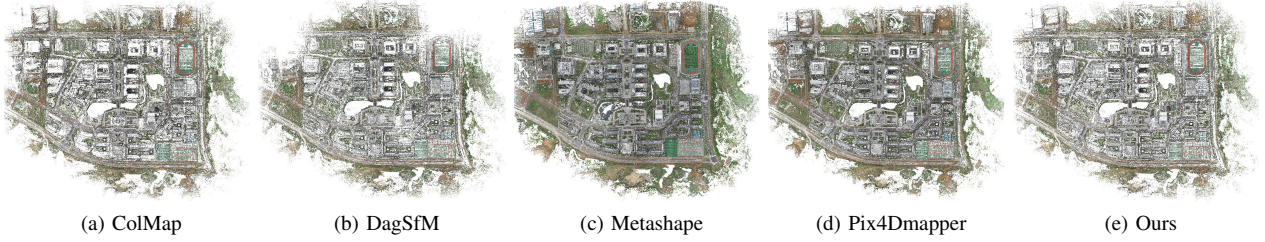(a) ColMap     (b) DagSfM     (c) Metashape     (d) Pix4Dmapper     (e) Ours

Fig. 10. The reconstructed sparse points of dataset 2 based on relative BA without GCPs.

TABLE III
THE STATISTICAL RESULTS OF ABSOLUTE BA WITH GCPS FOR DATASET 3.

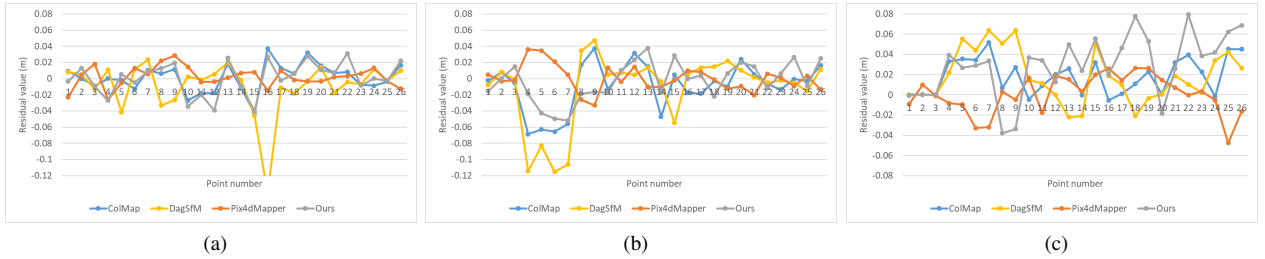| Method | Max (m) | | | Mean (m) | | | Std.dev. (m) | | |
|---|---|---|---|---|---|---|---|---|---|
| | $|X|$ | $|Y|$ | $|Z|$ | $|X|$ | $|Y|$ | $|Z|$ | $|X|$ | $|Y|$ | $|Z|$ |
| ColMap | 0.043 | 0.069 | 0.052 | 0.014 | 0.022 | 0.020 | 0.017 | 0.029 | 0.018 |
| DagSfM | 0.142 | 0.115 | 0.064 | 0.019 | 0.027 | 0.023 | 0.032 | 0.044 | 0.026 |
| Pix4dMapper | 0.028 | 0.036 | 0.048 | 0.010 | 0.012 | 0.015 | 0.013 | 0.016 | 0.019 |
| Ours | 0.040 | 0.052 | 0.080 | 0.016 | 0.020 | 0.036 | 0.020 | 0.024 | 0.031 |



(a)     (b)     (c)

Fig. 11. The individual residual plot in the X, Y, and Z directions for dataset 3: (a) the residual plot in the X direction; (b) the residual plot in the Y direction; (c) the residual plot in the Z direction.

in Figure 11a and Figure 11b. In addition, when compared with DagSfM, the proposed solution has achieved better accuracy in the horizontal direction for the three metrics. Considering that the GSD of dataset 3 is 1.2 cm, the overall geo-referencing accuracy of the proposed solution is approximately 1.7, 2.0, and 2.6 times the GSD value in the X, Y, and Z directions. Based on the results of absolute BA, we can conclude that

by using the WCDS-based global geometric constraint, the proposed parallel SfM solution can achieve more evenly distributed and comparative geo-referencing accuracy in the horizontal direction.

## V. CONCLUSIONS

In this paper, we have proposed the WCDS algorithm for the global model extraction and a parallel SfM solution to

implement efficient and accurate 3D reconstruction for UAV images. The core idea of the proposed solution is to create a global model to facilitate cluster merging. The experimental results from real UAV images demonstrate that the proposed parallel SfM can dramatically increase the reconstruction efficiency when compared with both open-source and commercial software packages, and comparative orientation accuracy has also been achieved in both relative and absolute BA tests.

Some limitations and possible improvements have also been observed in this study. First, the proposed SfM solution depends on the global model to assist scene merging. In the experiments as presented in Section IV-B, there are approximately 20% of images retained to build the global model. With the increase of datasets, the volume of the global model would become too large and then degenerates the overall performance of the parallel solution. To overcome the limitation, the global model would cooperate with the hierarchical strategy [34] to avoid increasing growth. Second, by the further analysis of time consumption, we found that approximately half of the time costs are consumed in the scene merging and the final BA optimization due to two main reasons. On the one hand, the bi-directional mean square reprojection error has been used in the RANSAC framework for the similarity transformation estimation, which is more time consuming than the direct point cloud-based methods; on the other hand, all camera poses and resumed 3D points participate in the final optimization, which causes high computation. To cope with this situation, future work would focus on designing an efficient strategy for similarity transformation estimation and selecting the most valuable tracks for the final BA optimization [23], [24].
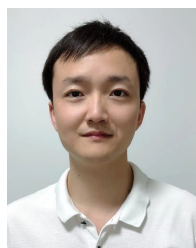
## ACKNOWLEDGMENT

## REFERENCES

[1] S. Jiang, W. Jiang, and L. Wang, "Unmanned aerial vehicle-based photogrammetric 3d mapping: A survey of techniques, applications, and challenges," *IEEE Geoscience and Remote Sensing Magazine*, 2021.

[2] S. Jiang, W. Jiang, W. Huang, and L. Yang, "Uav-based oblique photogrammetry for outdoor data acquisition and offsite visual inspection of transmission line," *Remote Sensing*, vol. 9, no. 3, p. 278, 2017.

[3] W. Huang, S. Jiang, and W. Jiang, "A model-driven method for pylon reconstruction from oblique uav images," *Sensors*, vol. 20, no. 3, p. 824, 2020.

[4] J. Zheng, H. Fu, W. Li, W. Wu, L. Yu, S. Yuan, W. Y. W. Tao, T. K. Pang, and K. D. Kanniah, "Growing status observation for oil palm trees using unmanned aerial vehicle (uav) images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 173, pp. 95–121, 2021.

[5] T. Bakirman, B. Bayram, B. Akpinar, M. F. Karabulut, O. C. Bayrak, A. Yigitoglu, and D. Z. Seker, "Implementation of ultra-light uav systems for cultural heritage documentation," *Journal of Cultural Heritage*, vol. 44, pp. 174–184, 2020.

[6] S. Jiang and W. Jiang, "Efficient sfm for oblique uav images: From match pair selection to geometrical verification," *Remote Sensing*, vol. 10, no. 8, p. 1246, 2018.

[7] N. Snavely, S. M. Seitz, and R. Szeliski, "Photo tourism: exploring photo collections in 3d," in *ACM transactions on graphics (TOG)*, vol. 25. ACM, 2006, pp. 835–846.

[8] H. Cui, X. Gao, S. Shen, and Z. Hu, "Hsfm: Hybrid structure-from-motion," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1212–1221.

[9] S. Jiang, C. Jiang, and W. Jiang, "Efficient structure from motion for large-scale uav images: A review and a comparison of sfm tools," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 167, pp. 230–251, 2020.

[10] C. Wu, S. Agarwal, B. Curless, and S. M. Seitz, "Multicore bundle adjustment," in *CVPR 2011*. IEEE, 2011, pp. 3057–3064.

[11] Y. Sun, H. Sun, L. Yan, S. Fan, and R. Chen, "Rba: Reduced bundle adjustment for oblique aerial photogrammetry," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 121, pp. 128–142, 2016.

[12] M. I. A. Lourakis, "Sba: A software package for generic sparse bundle adjustment," vol. 36, no. 1, p. 2, 2009.

[13] B. Bhowmick, S. Patra, A. Chatterjee, V. M. Govindu, and S. Banerjee, "Divide and conquer: Efficient large-scale structure from motion using graph partitioning," in *Asian Conference on Computer Vision*. Springer, 2014, pp. 273–287.

[14] N. Snavely, S. M. Seitz, and R. Szeliski, "Skeletal graphs for efficient structure from motion," in *2008 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2008, pp. 1–8.

[15] M. Havlena, A. Torii, and T. Pajdla, "Efficient structure from motion by graph optimization," *Computer vision–ECCV 2010*, pp. 100–113, 2010.

[16] S. Choudhary and P. Narayanan, "Visibility probability structure from sfm datasets and applications," in *European conference on computer vision*. Springer, 2012, pp. 130–143.

[17] R. Hänsch, I. Drude, and O. Hellwich, "Modern methods of bundle adjustment on the gpu." *ISPRS Annals of Photogrammetry, Remote Sensing & Spatial Information Sciences*, vol. 3, no. 3, 2016.

[18] X. Liu, W. Gao, and Z.-Y. Hu, "Hybrid parallel bundle adjustment for 3d scene reconstruction with massive points," *Journal of Computer Science and Technology*, vol. 27, no. 6, pp. 1269–1280, 2012.

[19] M. Zheng, S. Zhou, X. Xiong, and J. Zhu, "A new gpu bundle adjustment method for large-scale data," *Photogrammetric Engineering & Remote Sensing*, vol. 83, no. 9, pp. 633–641, 2017.

[20] E. Rupnik, F. Nex, and F. Remondino, "Automatic orientation of large blocks of oblique images," *International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 40, p. 1, 2013.

[21] A. Cefalu, N. Haala, and D. Fritsch, "Structureless bundle adjustment with self-calibration using accumulated constraints," *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial Information Sciences*, vol. 3, no. 3, pp. 3–9, 2016.

[22] A. L. Rodríguez, P. E. López-de Teruel, and A. Ruiz, "Reduced epipolar cost for accelerated incremental sfm," in *CVPR 2011*. IEEE, 2011, pp. 3097–3104.

[23] M. Cao, W. Jia, Z. Lv, Y. Li, W. Xie, L. Zheng, and X. Liu, "Fast and robust feature tracking for 3d reconstruction," *Optics and Laser Technology*, vol. 110, pp. 120–128, 2019.

[24] G. Zhang, H. Liu, Z. Dong, J. Jia, T.-T. Wong, and H. Bao, "Efficient non-consecutive feature tracking for robust structure-from-motion," *IEEE Transactions on Image Processing*, vol. 25, no. 12, pp. 5957–5970, 2016.

[25] R. Kummerle, G. Grisetti, H. Strasdat, K. Konolige, and W. Burgard, "G2o: A general framework for graph optimization," in *2011 IEEE International Conference onRobotics and Automation (ICRA)*, 2011.

[26] Y. Chen, S. Shen, Y. Chen, and G. Wang, "Graph-based parallel large scale structure from motion," *Pattern Recognition*, vol. 107, p. 107537, 2020.

[27] L. Lu, Y. Zhang, and K. Liu, "Block partitioning and merging for processing large-scale structure from motion problems in distributed manner," *IEEE Access*, vol. 7, pp. 114 400–114 413, 2019.

[28] B. Xu, L. Zhang, Y. Liu, H. Ai, B. Wang, Y. Sun, and Z. Fan, "Robust hierarchical structure from motion for large-scale unstructured image sets," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 181, pp. 367–384, 2021.

[29] S. Zhu, T. Shen, L. Zhou, R. Zhang, J. Wang, T. Fang, and L. Quan, "Parallel structure from motion from local increment to global averaging," *arXiv preprint arXiv:1702.08601*, 2017.

[30] J. Shi and J. Malik, "Normalized cuts and image segmentation," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 22, no. 8, pp. 888–905, 2000.

[31] H.-Y. Shum, Q. Ke, and Z. Zhang, "Efficient bundle adjustment with virtual key frames: A hierarchical approach to multi-frame structure from motion," in *Proceedings. 1999 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (Cat. No PR00149)*, vol. 2. IEEE, 1999, pp. 538–543.

[32] M. Farenzena, A. Fusiello, and R. Gherardi, "Structure-and-motion pipeline on a hierarchical cluster tree," in *2009 IEEE 12th International*

*Conference on Computer Vision Workshops, ICCV Workshops.* IEEE, 2009, pp. 1489–1496.

[33] R. Gherardi, M. Farenzena, and A. Fusiello, "Improving the efficiency of hierarchical structure-and-motion," in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition.* IEEE, 2010, pp. 1594–1600.

[34] R. Toldo, R. Gherardi, M. Farenzena, and A. Fusiello, "Hierarchical structure-and-motion recovery from uncalibrated images," *Computer Vision and Image Understanding*, vol. 140, pp. 127–143, 2015.

[35] K. Ni and F. Dellaert, "Hypersfm," in *2012 Second International Conference on 3D Imaging, Modeling, Processing, Visualization and Transmission.* IEEE, 2012, pp. 144–151.

[36] B. Bhowmick, S. Patra, A. Chatterjee, V. M. Govindu, and S. Banerjee, "Divide and conquer: A hierarchical approach to large-scale structure-from-motion," *Computer Vision and Image Understanding*, vol. 157, pp. 190–205, 2017.

[37] C. Wu, "Towards linear-time incremental structure from motion," in *2013 International Conference on 3D Vision-3DV 2013.* IEEE, 2013, pp. 127–134.

[38] R. Shah, A. Deshpande, and P. Narayanan, "Multistage sfm: Revisiting incremental structure from motion," in *2014 2nd International Conference on 3D Vision*, vol. 1. IEEE, 2014, pp. 417–424.

[39] J. Dong and S. Soatto, "Domain-size pooling in local descriptors: Dsp-sift," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 5097–5106.

[40] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, "Bundle adjustment: a modern synthesis," in *International workshop on vision algorithms.* Springer, 1999, pp. 298–372.

[41] S. Jiang and W. Jiang, "Efficient structure from motion for oblique uav images based on maximal spanning tree expansion," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 132, pp. 140–161, 2017.

[42] S. Jiang and W. Jiang, "On-board gnss/imu assisted feature extraction and matching for oblique uav images," *Remote Sensing*, vol. 9, no. 8, p. 813, 2017.

[43] S. Jiang, W. Jiang, and B. Guo, "Leveraging vocabulary tree for simultaneous match pair selection and guided feature matching of uav images," *ISPRS Journal of Photogrammetry and Remote Sensing*, vol. 187, pp. 273–293, 2022.

[44] D. Nister and H. Stewenius, "Scalable recognition with a vocabulary tree," in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*, vol. 2. IEEE, 2006, pp. 2161–2168.

[45] A. M. Andrew, "Another efficient algorithm for convex hulls in two dimensions," *Information Processing Letters*, vol. 9, no. 5, pp. 216–219, 1979.

[46] S. Guha and S. Khuller, "Approximation algorithms for connected dominating sets," *Algorithmica*, vol. 20, no. 4, pp. 374–387, 1998.

[47] S. Umeyama, "Least-squares estimation of transformation parameters between two point patterns," *IEEE Transactions on Pattern Analysis & Machine Intelligence*, vol. 13, no. 04, pp. 376–380, 1991.

[48] C. Wu, "Siftgpu: A gpu implementation of david lowe's scale invariant feature transform (sift)," https://github.com/pitzer/SiftGPU, 2007, accessed: 2017-06-19.

[49] J. L. Schonberger and J.-M. Frahm, "Structure-from-motion revisited," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4104–4113.

[50] Agisoft, "Agisoft metashape homepage," http://www.agisoft.com, 2021, accessed: 2021-12-18.

[51] Pix4Dmapper, "Pix4dmapper homepage," https://www.pix4d.com, 2021, accessed: 2021-12-18.

[52] X. Zhou, K. Xie, K. Huang, Y. Liu, Y. Zhou, M. Gong, and H. Huang, "Offsite aerial path planning for efficient urban scene reconstruction," *ACM Transactions on Graphics (TOG)*, vol. 39, no. 6, pp. 1–16, 2020.

**San Jiang** received the B.S. degree in remote sensing science and technology from Wuhan University in 2010, and the M.Sc. and Ph.D. degrees in photogrammetry and remote sensing from Wuhan Univeristy in 2012 and 2018, respectively. From 2012 to 2014, he worked as an assistant engineer in Tianjin Institute of Surveying and Mapping. From 2014 to 2015, he joined the LIESMARS (State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing of Wuhan Univeristy) as a research assistant.

Currently, he is an associate professor in the School of Computer Science at China University of GeoSciences (Wuhan). His research interests include image matching, SfM-based aerial triangulation, and 3D reconstruction.



**Qingquan Li** received the Ph.D. degree in Geographic Information System (GIS) and photogrammetry from the Wuhan Technical University of Surveying and Mapping, Wuhan, China, in 1998. From 1988 to 1996, he was an Assistant Professor with Wuhan University, Wuhan, where he became an Associate Professor, in 1996, and has been a Professor, since 1998. He is currently a Professor of Shenzhen University, Shenzhen, China; a Professor with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University; and the Director of Shenzhen Key Laboratory of Spatial Smart Sensing and Service, Shenzhen. His research interests include intelligent transportation systems, 3-D and dynamic data modeling, and pattern recognition.

Dr. Li is an Academician of International Academy of Sciences for Europe and Asia (IASEA), an Expert in Modern Traffic with the National 863 Plan, and an Editorial Board Member of the Surveying and Mapping Journal and the Wuhan University Journal—Information Science Edition.



**Wanshou Jiang** received his bachelor and master degrees in photogrammetry and remote sensing from Wuhan Technical University of Surveying and Mapping respectively in 1989 and 1996. In 2004, he received the PhD degree in photogrammetry and remote sensing from Wuhan University. He started his research career in 1989 as a software developer in analytical photogrammetry. In 2000, he joined the LIESMARS (the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing) as an associate researcher and then he got the tenure position of researcher in 2005.

His research interest includes image registering, image classification, change detection, 3D reconstruction, etc. He made a lot of contribution to the famous digital photogrammetric workstation VirtuoZo and designed a software platform, named OpenRS, for remote sensing image processing.



**Wu Chen** a Professor with the Department of Land Surveying and Geo-Informatics of the Hong Kong Polytechnic University. He has been actively working on GNSS related research for more than 30 years. His main research interests are geodesy and geodynamics, seamless positioning technologies, GNSS positioning and applications, system integration, GNSS performance evaluation, regional GPS network, and SLAM.