

# An Improved Augmented Reality Registration Method Based on Visual SLAM

Qing Hong Gao<sup>1</sup>, Tao Ruan Wan<sup>3</sup>, Wen Tang<sup>2(✉)</sup>, Long Chen<sup>2</sup>,  
and Kai Bing Zhang<sup>1</sup>

<sup>1</sup> College of Electronic and Information, Xian Polytechnic University, Xian, China

<sup>2</sup> Faculty of Science and Technology, Bournemouth University, Bournemouth, UK  
wtang@bournemouth.ac.uk

<sup>3</sup> Faculty of Informatics, University of Bradford, Bradford, UK  
t.wan@bradford.ac.uk

**Abstract.** Markerless Augmented Reality registration using standard Homography matrix is instable and has low registration accuracy. In this paper, we present a new method to improve the augmented reality registration method based on the Visual Simultaneous Localization and Mapping (VSLAM). We improved the method implemented in ORB-SLAM in order to increase stability and accuracy of AR registration. VSLAM algorithm generate 3D scene maps in dynamic camera tracking process. Hence, for AR based on VSLAM utilizes the 3D map of the scene reconstruction to compute the location for virtual object augmentation. In this paper, a Maximum Consistency with Minimum Distance and Robust Z-score (MCMD\_Z) algorithm is used to perform the planar detection of 3D maps, then the Singular Value Decomposition (SVD) and Lie group are used to calculate the rotation matrix that helps to solve the problem of virtual object orientation. Finally, the method integrates camera poses on the virtual object registration. We show experimental results to demonstrate the robustness and registration accuracy of the method for augmented reality applications.

**Keywords:** Augmented Reality · SLAM algorithm · Virtual registration and fusion · Point cloud

## 1 Introduction

Augmented Reality (AR) is the technology of mixing real scenes with virtual scenes, an emerging field of huge application potentials. The technology makes the use of computer-generated virtual information within the real world to enhance the human perception of the world. As defined by Azuma, it is an integration of virtual world and the real world with real-time interactions via three-dimensional registrations [2]. The recent rapid development of software as well as hardware technologies in virtual reality and computer vision, AR has a wider range of applications in medicine, military, entertainment and others [4, 15]. Virtual registration, however, remains a challenge in AR research. Initially

Simultaneous Localization and Mapping (SLAM), as a probability algorithm, has been mainly used for positioning robots in unknown environments [3, 14]. More recently, researchers have started to utilize the accuracy and real-time performance of SLAM for virtual registration in AR. Davison et al. [5, 6] have used a monocular camera to achieve fast 3D modeling and positioning of cameras in unknown environments, which has presented many practical uses of the algorithm. Klein [9] applied a SLAM algorithm in the creation of three-dimensional point clouds, as well as Reitmayr [12] demonstrated the use of SLAM and sensor fusion techniques in an accurate virtual reality registration with markerless tracking.

The method of computing homography matrix in AR systems for the three-dimensional registration [7, 11] is simple and efficient. This algorithm requires the detection of four point coordinates of a plane in order to determine the translation and rotation of the camera relative to the world coordinate system. In spite of its simplicity and efficiency, since it is based on the 2D plane registration, the four points of detection algorithm is prone to the error of misplacement of the virtual object registration, resulting in virtual objects being unstable with distracting visual effects (e.g. flashing visual artifacts). Previous approaches [9, 12] have attempted to make the use of the three-dimensional map information generated by SLAM for this process. In this paper, we present a method of improvement to the registration and tracking process of virtual objects by using map information generated by VSLAM [12] technology. The three-dimensional information of a scene generated by VSLAM cannot be used directly, due to the interference points and the large error of point clouds. Therefore, a robust Maximum Consistency with Minimum Distance and Robust Zscore (MCMD\_Z) [1] algorithm have been used to detect the 2D plane more accurately. Our improved MCMD\_Z method computes plane point matrix by using the plane normal vector fund by Singular Value Decomposition (SVD). A method of lie group is then used to convert the normal vector into the rotation matrix to register the virtual object using the plane information. We use the precise positioning function of the VSLAM to change the camera poses to the rendering coordinate system under the camera perspective for the three-dimensional registration of the virtual object.

The main contribution of this paper is to develop a method that can effectively produce stable and high registration accuracy for virtual reality fusion.

## 2 AR System Overview

The AR system consists of two software modules: VSLAM module and registration module as shown in Fig. 1 for an overview of the system. Tracking in the VSLAM module is to locate the camera position through processing each image frame, and decide when to insert a new keyframe. Firstly, the feature matching is initialized with the previous frame and Bundle Adjustment BA [16] is used to optimize the camera poses. While the 3D map is initialized and the map is successfully created by the VSLAM module, the registration module is called.

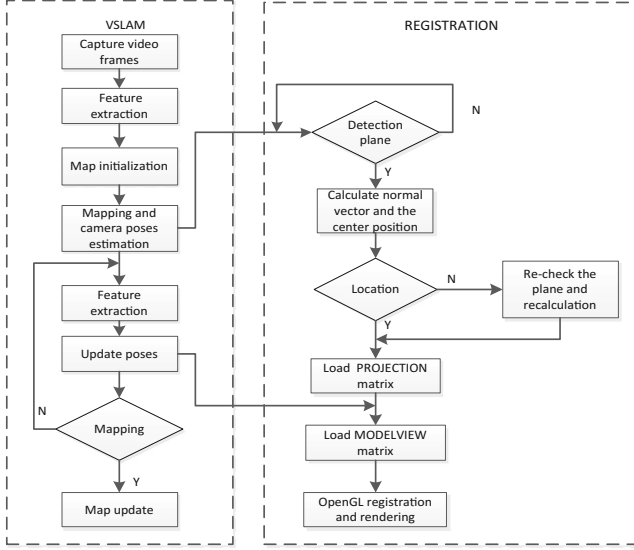


Fig. 1. System overview

Once the point cloud of the scene is generated, the computation for the plane detection is started, and the center and the normal vector of the plane are calculated. The center of the plane is used to determine the exact position of the virtual object and the normal vector is used to determine the orientation of the virtual object. Camera poses obtained by VSLAM are then transformed to the modelview matrix of OpenGL, which will be generated by the transformation of three-dimensional virtual objects to the center of the plane to achieve the virtual augmentation.

## 2.1 Tracking and 3D Map Building

Our system is based on a visual simultaneous mapping and tracking approach by extracting and matching the Oriented Features From Accelerated Segment Test (FAST) and Rotated Binary Robust Independent Elementary Features (BRIEF) (ORB) [13] feature points and compute two models: a homography matrix that is used to compute a planar scene and a fundamental matrix that is used to compute a non-planar scene. Each time two matrices are calculated, and a score ( $M = H$  for the homography matrix,  $M = F$  for the fundamental matrix) is also calculated as shown in Eq. 1. The score is used to determine which model is more suitable for the camera posture.

$$S_M = \sum_i (\rho_M (d_{cr,M}^2 (x_c^i, x_r^i)) + \rho_M (d_{rc,M}^2 (x_c^i, x_r^i))) \quad (1)$$

$$\rho_M (d^2) = \begin{cases} \Gamma - d^2 & \text{if } d^2 < T_M \\ 0 & \text{if } d^2 \geq T_M \end{cases}$$

where  $d_{rc}$  and  $d_{cr}$  is the measure of symmetric transfer errors [8],  $T_m$  is the outlier rejection threshold based on the  $\chi^2$ ,  $\Gamma$  is equal to,  $T_m$ ,  $x_c$  is the features of the current frame,  $x_r$  is the features of the reference frame. The BA is used to optimize camera poses, which gets a more accurate camera position as in following equation.

$$\{R, r\} = \arg \min_{R, t} \sum_{i \in \mathcal{X}} \rho \left( \|x^i - \pi(RX^i + t)\|_{\Sigma}^2 \right) \quad (2)$$

where  $R \in \mathcal{SO}^3$  is the rotation matrix,  $t \in \mathbb{R}^3$  is the translation vector,  $\xi^i \in \mathbb{R}^3$  is a three-dimensional point in space,  $x^i \in \mathbb{R}^2$  is the key point,  $\rho$  is the Huber cost function, Sigma item is the covariance matrix associated to the key point,  $\pi$  is the projection function.

After obtaining the accurate position estimation of the camera, the three-dimensional map point cloud is obtained by triangulating the key frame through the camera poses, and finally the local BA is used to optimize the map. A detailed description of the approach is given in [10].

### 3 Plane Detection and Calculation of the Normal Vector

The map created in §2.1 is composed of a sparse point cloud. Because of the error of the point cloud data with large number of abnormal values, MCMD\_Z is used for plane detection. A MCMD\_Z algorithm is used to fit the data according to a search model. The idea of this algorithm is to use Principal Component Analysis (PCA) for a reliable selection of the registration plane, using Robust Z-score to remove invalid points at once. This method not only effectively avoids the threshold setting, but also runs fast. The MCMD\_Z algorithm detects the plane as follows:

---

#### Algorithm 1. The MCMD\_Z algorithm

---

- 1: loop:
  - 2: Randomly select any 3 points in the original point cloud, and calculate the corresponding initial plane and its *normal vector*;
  - 3: Ranking the points according to the value of  $d_i$  that are the distance of the map point cloud to the initial plane;
  - 4: Select the threshold  $t = 2\sigma$ .  $\sigma$  is the standard deviation of the point cloud to the initial plane model distance. When the  $d_i > t$ , the point is removed as an exception point, and vice versa as valid data retention, the number of valid data,  $P_{num}$ ;
  - 5: **if** The set of points reach the minimum Eigen values **then**
  - 6:   Stop the loop.
  - 7: Calculate Robust Z-score uses the Eq. 3. Zscore greater than 2 is considered an outlier and will be removed.
-

$$Rz_i = \frac{\left| od_i - \underset{j}{\text{median}}(od_i) \right|}{a \cdot \underset{i}{\text{median}} \left| od_i - \underset{j}{\text{median}}(od_i) \right|} \quad (3)$$

The detection of the plane determines the plane location, while providing a super-position of the location for a virtual object. Although the location of the virtual object is determined, virtual objects will not appear parallel to the plane but at a certain angle to the plane. In order to solve this problem, we need to calculate the normal vector of the plane and the rotation matrix.

The SVD of the matrix in the plane interior point is obtained, and the right singular vector corresponding to the minimum eigenvalue is the normal vector of the plane. Because there are two normal vectors in the plane, it is important that the normal vector direction is pointing outward. Specifically, the vector of the camera to the plane is found by the camera's posture. Through the vector and the relationship between the plane vectors, we can then determine the direction of the normal vector. The rotation matrix is obtained from the known normal vector by Lie group using the following equation:

$$R_{3 \times 3} = \exp(\hat{w}) = I + \sin(\|w\|) \cdot \frac{\hat{w}}{\|w\|} + (1 - \cos(\|w\|)) \cdot \frac{\hat{w}\hat{w}}{\|w\|^2} \quad (4)$$

$$w = \frac{n_y \times n_p}{\|n_y \times n_p\|} \cdot \arctan \frac{\|n_y \times n_p\|}{n_y \times n_p}$$

Where  $n_y$  is  $y$ -axis unit vector,  $n_p$  is normal vector of the plane,  $w$  is a column vector,  $\hat{w}$  is the anti-symmetric matrix of the vector  $w$ . Finally, the transformation matrix of OpenGL is composed of a translation vector and a rotation matrix. The rotation matrix is obtained and the translation vector is found to be the center of the plane.

### 3.1 Virtual Registration

The virtual object is finally registered in the real world, which must go through the transformation of the coordinate systems (from the world coordinate system to the camera coordinate system to the crop coordinate system, and to the screen coordinate system). The transformation sequences can be described by Eq. 4 from left to right: the world coordinate system is transformed into the camera coordinate system by a rotation matrix  $R_{(3 \times 3)}$  and a translation matrix  $T_{(3 \times 1)}$ . Those matrices are made up by the camera's position and the detected plane information. Then the camera coordinate system is then transformed into the screen coordinate system  $(u, v)$  by the focal length  $(f_x, f_y)$  and the principal point  $(d_x, d_y)$ . These parameters are obtained by the camera calibration. Finally, the virtual object is registered in the screen to the real world.

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} f_x & 0 & d_x & 0 \\ 0 & f_y & d_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} R_{3 \times 3} & T_{3 \times 1} \\ 0_{1 \times 3} & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} \quad (5)$$

## 4 Experiment and Evaluation

Our experiment is run under Ubuntu 14.04 system, CPU clocked at 2.3 GHz, 8 GB memory and graphics card for the NVIDIA GeForce GTX 960 MB. The camera resolution is 640 by 480 pixels at 30 Hz. The experimental scene is indoors and the length of the image collection is 1857 frames. Figure 2(a)–(b) show the indoor scene under the AR tracking and registration. We can see that the tracking and registration effect. Figure 2(c) shows the correct virtual object orientation.



**Fig. 2.** AR tracking and registration (left to right (a)–(c))

### 4.1 Plane Detection Analysis

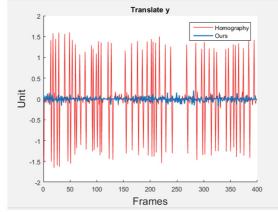
Our method based on the CMCD.Z, which achieves better results than Random Sample Consensus (RANSAC). In contrast to these two algorithms, we use the Gaussian distribution to produce 1000 point with outlier percentages (10 and 20) using the same input parameters used previously. The inliers have means (15.0, 15.0, 10.0) and variances (10.0, 2.0, 0.5). The outliers have means (15.0, 15.0, 10.0) and variances (10.0, 2.0, 0.5). The program ran 1000 times. We compared Correct Identification Rate (CIR) and Swamping Rate (SR). The RANSAC sets iterations 50 times (Table 1).

### 4.2 Registration Error Analysis

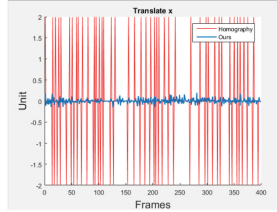
A comparison method is used with fixed camera positions to evaluate the robustness of the method. The three-dimensional registration of the virtual object is carried out by using the described method and the standard homography matrix method. Six components of the three-dimensional registration results are analyzed. The difference between the transformation matrix of the current frame

**Table 1.** Correct Identification Rate (CIR), Swamping Rate (SR) and Time

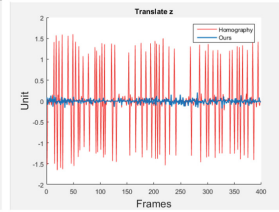
Methods	Outlier percentage				Time (s)
	10		20		
	CIR	SR	CIR	SR	
RANSAC	100	24.3	100	8	0.307849
MCMD_Z	100	6	100	0.5	0.116869



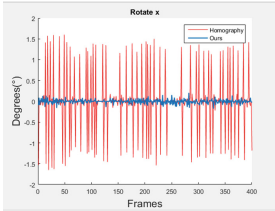
(a)



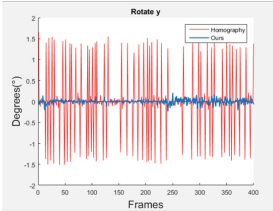
(b)



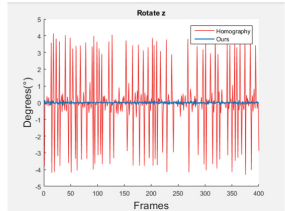
(c)



(d)



(e)



(f)

**Fig. 3.** Registration error (Color figure online)

and the corresponding component of the transformation matrix of the previous frame is used as the basis for comparison. The results are shown in Fig. 3, where Translate x, Translate y and Translate z are the errors of the translation components, respectively, and Rotate x, Rotate y, Rotate z are relative to the x, y, z axis of the rotation component error and where is obtained by subtracting the previous frame from the current frame. The result of the rotation component is obtained by dividing the respective components with the dot product of

the corresponding coordinate axis, and the translation component is the result obtained by the normalization process.

In Fig. 3, the red curves in the figures are the results of using only the homography matrix. The blue curves are the results of the new registration method used in this paper. As can be seen from the Fig. 3, the use of the homography matrix method to register the virtual objects has produced large fluctuations of registration errors that are equivalent to virtual object registration instability. However, the new method tested on each rotation component has been kept the error in a small range below  $0.5^\circ$ . The errors with Translate x, Translate y and Translate z are also small similar to the result of the rotation components.

Through the experimental results, it can be seen that the new method produces stable virtual registration and solve the flickering phenomenon in the virtual reality registration, hence, improves the stability of the AR system.

## 5 Conclusions and Future Work

This paper presents a stable and realistic tracking method based on three-dimensional map information generated by VSLAM method to track the registration of virtual objects to ensure the stability and real-time performance of registration. Our proposed method is faster and is able to achieve more accurate registration results. The experimental results show that the proposed method can effectively suppress the virtual object jitter, have a higher tracking accuracy with good tracking performance. The current three-dimensional map used in this paper is a sparse point cloud, which can only access limited space configuration information.

While this work has served to propose and prototype with experiments to show the effectiveness of the proposed approach, future work will consider the use of dense point cloud based on our proposed method.

**Acknowledgement.** This work is supported by Shanxi Province Science and Technology Department: Projects (2016JZ026 2016KW-043) and (2016GY-047).

## References

1. Azuma, R., Baillot, Y., Behringer, R., Feiner, S., Julier, S., MacIntyre, B.: Recent advances in augmented reality. *IEEE Comput. Graphics Appl.* **21**(6), 34–47 (2001)
2. Azuma, R.T.: A survey of augmented reality. *Presence Teleoperators Virtual Environ.* **6**(4), 355–385 (1997)
3. Bailey, T., Durrant-Whyte, H.: Simultaneous localization and mapping (SLAM): part II. *IEEE Robot. Autom. Mag.* **13**(3), 108–117 (2006)
4. Bimber, O., Raskar, R.: *Spatial Augmented Reality Merging Real and Virtual Worlds*. A K Peters Ltd., Natick (2005)
5. Davison, A.J., Mayol, W.W., Murray, D.W.: Real-time localization and mapping with wearable active vision. In: *Proceedings of Second IEEE and ACM International Symposium Mixed and Augmented Reality*, pp. 18–27, October 2003



6. Davison, A.J., Mayol, W.W., Murray, D.W.: Real-time workspace localisation and mapping for wearable robot. In: Proceedings of Second IEEE and ACM International Symposium Mixed and Augmented Reality, pp. 315–316, October 2003
7. Fiala, M.: Artag, a fiducial marker system using digital techniques. In: Proceedings of IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR 2005), vol. 2, pp. 590–596, June 2005
8. Hartley, R., Zisserman, A.: *Multiple View Geometry in Computer Vision*. Cambridge University Press, Cambridge (2000)
9. Klein, G., Murray, D.: Parallel tracking and mapping for small AR workspaces. In: Proceedings of the 6th IEEE and ACM International Symposium on Mixed and Augmented Reality, pp. 225–234, November 2007
10. Mur-Artal, R., Montiel, J.M.M., Tards, J.D.: Orb-SLAM: a versatile and accurate monocular SLAM system. *IEEE Trans. Robot.* **31**(5), 1147–1163 (2015)
11. Prince, S.J.D., Xu, K., Cheok, A.D.: Augmented reality camera tracking with homographies. *IEEE Comput. Graphics Appl.* **22**(6), 39–45 (2002)
12. Reitmayr, G., Eade, E., Drummond, T.W.: Semi-automatic annotations in unknown environments. In: Proceedings of the 6th IEEE and ACM International Symposium on Mixed and Augmented Reality, pp. 67–70, November 2007
13. Rublee, E., Rabaud, V., Konolige, K., Bradski, G.: Orb: an efficient alternative to sift or surf. In: Proceedings of International Conference on Computer Vision, pp. 2564–2571, November 2011
14. Strasdat, H., Montiel, J., Davison, A.J.: Visual SLAM: why filter? *Image Vis. Comput.* **30**(2), 65–77 (2012)
15. Szalaviri, Z., Gervautz, M.: The personal interaction panel - a two-handed interface for augmented reality. *Comput. Graph. Forum* **16**(3), C335–C346 (1997). <http://onlinelibrary.wiley.com/doi/10.1111/1467-8659.00137/abstract>
16. Triggs, B., McLauchlan, P.F., Hartley, R.I., Fitzgibbon, A.W.: Bundle adjustment—a modern synthesis. In: Triggs, B., Zisserman, A., Szeliski, R. (eds.) *IWVA 1999*. LNCS, vol. 1883, pp. 298–372. Springer, Heidelberg (2000). doi:[10.1007/3-540-44480-7\\_21](https://doi.org/10.1007/3-540-44480-7_21)