

Terrain Aided Planetary UAV Localization Based on Geo-referencing

Xue Wan[✉], Yuanbin Shao[✉], Shengyang Zhang, and Shengyang Li[✉]

Abstract—The autonomous real-time optical navigation of planetary unmanned aerial vehicle (UAV) is of the key technologies to ensure the success of the exploration. In such a GPS-denied environment, vision-based localization is an optimal approach. In this article, we proposed a terrain aided simultaneous localisation and mapping (SLAM) algorithm, which simultaneously reconstructs the 3-D map point of environment and estimates the location of a planet UAV based on preexisting digital elevation model (DEM). To directly georeference the onboard UAV images to the digital terrain model, a theoretical model is proposed to prove that topographic features of UAV image and DEM can be correlated in the frequency domain via cross power spectrum. To provide the six-DOF of the UAV, we developed an optimization approach, which fuses the geo-referencing result into an SLAM system via local bundle adjustment (LBA) to achieve robust and accurate vision-based navigation even in featureless planetary areas. To test the robustness and effectiveness of the proposed localization algorithm, a new dataset for planetary drone navigation is proposed based on simulation engine. The proposed dataset includes 40 200 synthetic drone images taken from nine planetary scenes with related DEM query images. Comparison experiments are carried out to demonstrate that over the flight distance of 33.8 km, the proposed method achieved an average localization error of 0.45 m, compared to 1.32 m by ORB-SLAM2 and 0.75 m by ORB-SLAM3, with the processing speed of 12 Hz, which will ensure real-time performance. We will make our datasets available to encourage further work on this topic.

Index Terms—Frequency, multimodal registration, simultaneous localisation and mapping (SLAM), unmanned aerial vehicle (UAV) localization.

NOMENCLATURE

- x Horizontal coordinate in the image space domain.
- y Vertical coordinate in the image space domain.
- a Amount of translation in the horizontal direction.
- b Amount of translation in the vertical direction.
- u Horizontal coordinate in the image frequency domain.
- v Vertical coordinate in the image frequency domain.

Manuscript received 17 February 2022; revised 16 June 2022; accepted 27 July 2022. Date of publication 16 August 2022; date of current version 30 August 2022. This work was supported by the National Natural Science Foundation of China (NSFC) under Grant 4217012350. (Corresponding author: Yuanbin Shao.)

Xue Wan and Shengyang Li are with the Key Laboratory of Space Utilization and the Technology and Engineering Center for Space Utilization, Chinese Academy of Sciences, Beijing 100094, China (e-mail: wanxue@csu.ac.cn; shyy@csu.ac.cn).

Yuanbin Shao and Shengyang Zhang are with the School of Aeronautics and Astronautics, University of Chinese Academy of Sciences, Beijing 101408, China (e-mail: shaoyuanbin20@mails.ucas.ac.cn; zsy13947168513@163.com).

Digital Object Identifier 10.1109/TGRS.2022.3198745

ω	Radius coordinates in polar coordinates.
θ	Angular coordinates in polar coordinates.
τ	Azimuth angle of solar radiation.
σ	Elevation angle of solar radiation.
r	Reflectance of terrain.
δ	Function of peak value.
IFT	Function of inverse Fourier transform.
I_i	i th image.
F_i	Fourier transforms of the i th image.
Q	Cross power spectrum of image.
V_H	Elevation model.
F_H	Fourier spectrum of elevation model.
F_H^*	Complex conjugate of F_H .
R	3-D rotation matrix.
t	3-D translation vector.
T	3-D transformation matrix.
m_k^i	Coordinates of the i th 2-D feature points on the image I_k .
M_j	3-D position of the j th 3-D map point.
P_k	Camera 3-D pose at time k .
t_{vo}	Camera trajectory in the local reference frame.
S_k^t	True 3-D position in the local reference frame at time k .
t_k	3-D position transformation of camera at time k .
I_k	Unmanned aerial vehicle (UAV) image frame at time k .
R_k	Reference terrain shading image at time k .
τ_k	Azimuth of the UAV at time k .
T_{pl}	Planned flight path of UAV.
S_k^R	Reference 3-D position of the camera at time k .
T_{vo}	Camera trajectory in the global reference frame.
S_k^T	True 3-D position in the global reference frame at time k .
H	Height of the UAV to the ground surface.
w	Width of the UAV image.
h	Height of the UAV image.
H_{Tt}	Transformation between the local reference frame and the global reference frame.
$e_{i,j}$	Reprojection error of 3-D point M_j on the image I_i .
w_i	Confidence of observation m_k^i .

e_i^G	Deviations between geo-referencing and visual odometry.
w_i^G	Confidence of geo-referencing.
$\pi(\cdot)$	Projection function from 3-D scene to 2-D image.
$SE(3)$	3-D special Euclidean group.
$PC(\cdot)$	Phase correlation function.
GSD	Ground sample distance.
$DEM(\cdot)$	Terrain evaluation function.

I. INTRODUCTION

PLANETARY exploration is important to the answers to many fundamental questions, including the formation of the universe, the evolution of life, and the origin of the Earth. Unmanned aerial vehicles (UAVs) can reach places and areas where man is unable or very difficult to access. For planetary exploration, UAVs can be an effective tool for investigation and survey. The National Aeronautics and Space Administration (NASA) has led a project titled “BEES for Mars,” which utilizes UAVs for Mars exploration [1]. Ingenuity, the Mars helicopter [2] has succeeded in landing on Mars surface together with Perseverance rover and flied for 18 times already since December 31, 2021. NASA also planned to send a mobile robotic rotorcraft named Dragonfly to Titan, the largest moon of Saturn [3].

As one of the key technologies in planetary exploration, robust and accurate navigation systems ensure the safety and precision of scientific tasks. Unlike on Earth, GPS service is unavailable for planetary exploration, and thus, optical-based navigation has become one of the key technologies in autonomous navigation [4]. There are several benefits in optical-based navigation, including low power consumption, lightweight, and small size. Optical sensors can be used for multiple purposes, including obstacle avoidance [5], target area detection [4], and 3-D terrain surface reconstruction [6].

As there is hardly any artificial architecture on the planetary surface, terrain feature becomes one of the fundamental and crucial processing units in vision-based navigation on the planetary surface. To utilize terrain features for planetary navigation, several approaches have been proposed, including skyline matching [7], [8], ray tracing [9], and 3-D topography matching [10]. However, none of these methods can solve the multimodal matching problem between topography and optical imagery, and in this article, we aimed to answer the question: “**How can we find the shared features between digital terrain model and optical UAV imagery?**”

Directly matching the optical images and digital elevation model (DEM) data has the following merits. First, DEM data can be acquired from stereo-matching or SAR interference and is always available for planetary exploration. Second, as terrain feature is important in planetary navigation, DEM data are the representation of the terrain model, so terrain features can be directly extracted and matched. Third, the illumination effect between UAV and satellite can be largely eliminated. The reference map can be generated under a similar illumination condition as UAV images.

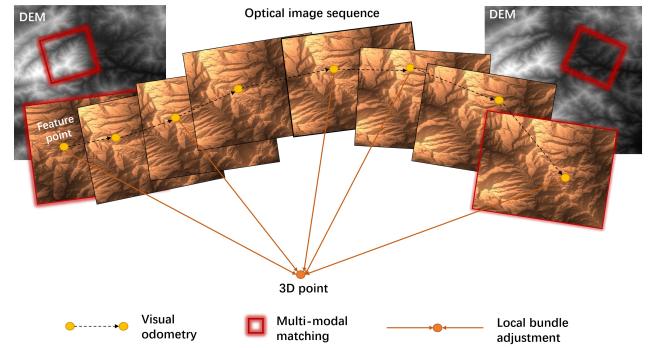


Fig. 1. Geo-referencing-based optical navigation for planetary UAV navigation.

In this article, we propose a terrain aided navigation approach based on frequency domain phase correlation (PC), which can localize the UAV in a terrain elevation map. The global localization results will be fused with visual odometry estimation results via local bundle adjustment (LBA) by proposing a new cost function, as shown in Fig. 1.

The contributions of this article are summarized as follows.

1) This article proposed a new cross-source planetary UAV localization dataset, including 40 200 synthetic drone images with ground-truth terrain models taken under nine different planetary topography scenes.

2) We proved via mathematical derivation that the Fourier spectrum of optical images can be divided into illumination and topography features, and the topography features of the terrain model and optical imagery can be correlated via frequency-based correlation.

3) This article proposes a **terrain aided optical navigation algorithm, terrain-simultaneous localisation and mapping (SLAM)**, which combined odometry estimation and frequency-based geo-referencing via a newly designed LBA framework.

II. RELATED WORK

One of the key challenges in planetary UAV optical navigation is to find a suitable image matching/image retrieval method for multimodal planetary images. All the approaches can be divided into five categories: **feature-based matching, area-based matching, skyline matching, 2-D–3-D ray tracing, and 3-D topography matching**.

Feature matching is one of the most commonly used matching algorithms for planetary optical navigation. Fixed-point (FP) feature [11] and Shi–Tomasi feature-based image matching [12] has been applied in Itokawa asteroid exploration. Harris corner detection and feature matching are used for Mars rover, Curiosity, and Opportunity, in the landing process [13]. SIFT feature matching has been applied in the CE-3 landing on the moon [14]. Speeded up robust features (SURFs) [15] improve feature extraction and description in SIFT to become more efficient. With the increasing real-time requirements of vision-based navigation, the application of faster oriented FAST and rotated BRIEF (ORB) features has become widespread and applied in ORB-SLAM2 [16] and ORB-SLAM3 [17]. Recently, deep-learning-based matching,

TABLE I
SUMMARY OF STATE-OF-THE-ART PLANETARY NAVIGATION APPROACHES

Methods	Perspective robustness	Scale robustness	Illumination robustness	Multi-modal robustness	Computational efficient
Feature-based	✓	✓	✗	✗	✗
Area-based	✗	✓	✗	✗	✗
Skyline matching	✓	✓	✗	✗	✓
Ray tracing	✓	✗	✓	✗	✓
3D topography matching	✓	✗	✓	✗	✗

such as SuperPoint [18], has been introduced in image matching. Characteristics, such as self-supervised training, homographic adaptation, and cross-domain adaptation in SuperPoint, have made great progress in terms of efficiency and accuracy in visual SLAM for robots. Luigi Freda implemented SURF-SLAM and SuperPoint-SLAM based on SURF and SuperPoint features in his proposed pySLAM,¹ respectively. Feature-based matching, which has been improved for years in the computer vision society, can cope with relatively large geometric image distortions between images. However, they can hardly cope with featureless areas that commonly exist on the planetary surface.

Another type of matching algorithm is area- or correlation-based algorithms. Typical area-based matching algorithms include normalized cross correlation (NCC), mutual information (MI) [19], and local self-similarities (LSSs) [20]. Sum of absolute difference (SAD) has been used for Mars rover obstacle avoidance and further improved by Chilian and Hirschmüller [21] using five matching windows for the task of Mars terrain 3-D reconstruction. Compared to feature-based algorithms, area-based matching algorithms are more robust to the featureless area as they take global gray value distribution. However, they are sensitive to camera pose variation, which makes them mostly applied in stereo vision instead of optical navigation in planetary exploration.

Another type of solution is to use topography features, as they are distinctive and unique and largely exist even in featureless areas, and thus, DEM is widely used as a reference map for absolute localization for planetary exploration. The most common DEM-aided navigation is skyline matching [7], [8]. These approaches first extract a skyline from the rover image, and then, a simulated skyline is rendered at an estimated rover position in the DEM. Finally, similarity values between the extracted skylines and the simulated skylines are calculated. The high similarity value stands for the high likelihood of the estimated rover position. Although skyline matching is commonly used to localize a rover within a known DEM, it is not suitable for planetary UAV navigation as the view angle of UAV is nadir and the skyline may not be visible. Moreover, the accuracy of skyline matching is largely determined by the DEM resolution, which may lead to localization errors if the resolution of DEM is not high enough.

Another type of terrain relative navigation approach is 3-D topography registration between the DEM and the

generated point cloud from onboard sensors, such as cameras and LiDAR. While dense 3-D point cloud can be directly generated from LiDAR scanning, the high-energy requirement prevents it from wide utilization in planetary navigation. The 3-D point cloud can also be generated from stereo matching, and however, they may suffer from the featureless areas. Even though the 3-D point cloud can be successfully generated, how to register it with DEM is another problem. One of the commonly used 3-D point cloud registration algorithms is iterative closest point (ICP) [10], [22], but ICP can hardly work with occlusion when a rover is facing a small hill. To overcome the limitation of ICP, Fang [9] proposed a ray-tracing approach to register stereo imagery to a lunar terrain model. This approach, however, still suffers from variations between the generated 3-D point cloud and DEM due to the largely different viewing distances and angles.

The merits and disadvantages of the five state-of-the-art planetary navigation approaches are summarized in Table I. As shown in Table I, none of the approaches can cope with the multimodal difference between the terrain model and optical UAV images.

In this article, we propose to use frequency information for multimodal matching to localize the current view of UAV in a digital elevation map. To the best of our knowledge, none of the above approaches try to directly match the optical image with the DEM. This may be because the features extracted from optical images and DEM are so different that it becomes really difficult to correlate them. In this article, we demonstrated that although it is visually difficult to locate the optical image in DEM, their frequency information can be correlated via PC, and thus, the multimodal image matching can be done in the frequency domain.

Then, the next question goes to how to iteratively estimate the six DOF of the planetary UAV. In this article, we seamlessly incorporate multimodal registration into an LBA algorithm by proposing to use the phase information to correlate the topography features between DEM data and optical images. Finally, we proposed that a terrain aided SLAM is proposed by seamlessly integrating the multimodal registration result into a visual odometry estimation.

III. THEORETICAL ANALYSIS FOR FREQUENCY-BASED TERRAIN FEATURE MATCHING

A. Phase Correlation

This section will briefly introduce the principle and history of PC, a frequency-based image matching algorithm that estimates the image shift between two images via the shift

¹<https://github.com/luigifreda/pyslam>

property of Fourier transform (FT). According to the shift property of FT, the translation shift (a, b) in the image spatial domain results in a linear phase shift in the Fourier frequency domain. Then, PC is defined as the normalized cross power spectrum between the FTs of the two images $I_1(x, y)$ and $I_2(x, y)$

$$Q(\omega, \theta) = \frac{F_H^*(\omega, \theta)}{Z_d} e^{-i(au+bv)} \quad (1)$$

where $*$ stands for complex conjugate. $F_1(u, v)$ and $F_2(u, v)$ in (u, v) coordinates are the FTs of the two images $I_1(x, y)$ and $I_2(x, y)$, respectively.

The shifts (a, b) can be resolved at integer level via inverse Fourier transform (IFT) to convert $Q(u, v)$ into an approximate Dirac delta function as

$$\text{IFT}(Q(u, v)) = \delta(x - a, y - b). \quad (2)$$

The translation (a, b) can also be solved directly in the frequency domain with subpixel accuracy by unwrapping and fitting the PC fringes in the cross power spectrum $Q(u, v)$ [23], [24], [25] or spatial domain by fitting the peak of function δ to a Gaussian function [26]. The peak value of function δ , ranging from 0 to 1, indicates the quality of PC matching. If the two images are identical, the peak value of δ equals 1.

When the rotation between $I_1(x, y)$ and $I_2(x, y)$ is θ_0 , the rotation θ_0 can be resolved at an integer level via transformation from Cartesian coordinates to polar coordinates and IFT to convert $Q(u, v)$ into an approximate Dirac delta function as

$$\text{IFT}(Q(u, v)) = \delta(\rho, \theta - \theta_0). \quad (3)$$

In our previous work, we have proved the robustness of PC to sun angle variation and shadow via the PC decomposition (PCD) model [27]. Using the illumination-invariant property of PC, we applied it to illumination-invariant change detection based on pixel-to-pixel matching [28] and UAV navigation based on the geo-referencing between UAV image and reference satellite images [29]. To enhance the robustness of PC toward the complex geometric distortion, we combine PC with particle swarm optimization [30]. We demonstrate the superior performance of PC toward state-of-the-art image matching algorithms in UAV navigation tasks [30]. In this article, we will further investigate the robustness of PC to multimodal data, including DEM and optical imagery via theoretical analysis. While our previous work focused on the image matching algorithm, this article will propose a fully terrain aided navigation pipeline, which combines geo-referencing results and visual odometry results for planetary UAV navigation.

B. Robustness of PC for DEM-Based Geo-referencing

In this section, we will dig out the relationship between UAV optical image and the topography model in the frequency domain and then prove via mathematical derivation that the terrain feature shared by the optical image and the topography model can be extracted and correlated via PC.

First, we will figure out the mathematical relationship between the optical image UAV image $I(x, y)$ and the topography model $V_H(x, y)$ in the spatial domain. As the planetary surface is a natural landscape, in this article, it is regarded as a Lambertian model. Under a given solar radiation intensity L in the direction of azimuth angle τ and elevation angle σ , the intensity of UAV image $I(x, y)$ generated from the elevation model is $V_H(x, y)$ [31]

$$I(x, y) = L \frac{p \cos \tau \cos \sigma + q \sin \tau \cos \sigma + \sin \sigma}{\sqrt{p^2 + q^2 + 1}} r(x, y) \quad (4)$$

where $p = (\partial/\partial x)V_H(x, y)$ and $q = (\partial/\partial y)V_H(x, y)$ are the gradients of $V_H(x, y)$ in the x - and y -directions and $r(x, y)$ represents the reflectance value at position (x, y) .

To analyze their relationship in the frequency domain, the UAV image $I(x, y)$ is transformed into the frequency domain as

$$F_I(\omega, \theta) = L[\sin \sigma + \cos \sigma \cos(\theta - \sigma)F_H(\omega, \theta)] \quad (5)$$

where $F_I(\omega, \theta)$ is the Fourier spectrum of elevation $I(x, y)$ and $F_H(\omega, \theta)$ is the Fourier spectrum of elevation $V_H(x, y)$.

It should be noted that here, we use a polar coordinate instead of a Cartesian coordinate to better illustrate the illumination effect on the frequency domain.

Then, the PC cross power spectrum of optical image $F_I(\omega, \theta)$ and DEM $F_H(\omega, \theta)$ is

$$Q(\omega, \theta) = \frac{\cos \sigma}{Z_d} \cos(\theta - \tau) + \frac{\sin \sigma}{Z_d} F_H^*(\omega, \theta) \quad (6)$$

where $Z_d = |\cos \sigma F_H^*(\omega, \theta) + \sin \sigma \cos(\theta - \sigma)|$.

As shown in (6), the cross power spectrum of a UAV optical image and reference DEM $Q(\omega, \theta)$ is mainly altered by two terms: the illumination term, $\cos(\theta - \sigma)$, and the topography term, $F_H^*(\omega, \theta)$. The following will investigate the effect of each term.

1) *Effect of Illumination Term:* To analyze the effect of illumination term, $\cos(\theta - \tau)$, on DEM-based geo-referencing, we generated a reference image of a DEM shown in Fig. 2(a) and a target image of a terrain shading image generated from the same DEM data under a small elevation angle $\sigma = 10^\circ$. In this case, term $\cos \sigma \cos(\theta - \sigma)$ becomes dominant in the PC cross power spectrum $Q(\omega, \theta)$ with other terms that are nearly zero, and then, $Q(\omega, \theta)$, shown in Fig. 2(c), can be simplified as

$$Q(\omega, \theta) = \frac{\cos \sigma}{Z_d} \cos(\theta - \tau). \quad (7)$$

Then, we can consider that there is a shift between $I(x, y)$ and $V_H(x, y)$, so their cross power spectrum, shown in Fig. 2(d), becomes

$$Q(\omega, \theta) = \frac{\cos \sigma}{Z_d} \cos(\theta - \tau) e^{-i(au+bv)}. \quad (8)$$

As shown in Fig. 2(d), several fringe patterns are flipped due to the term $\cos(\theta - \sigma)$. If the value of $\cos(\theta - \sigma)$ is negative, the fringes will be flipped. However, for both positive and negative correlation in the PC, the cross power spectrum remains the same, and thus, the fringe density $(a^2 + b^2)^{1/2}$ and orientation b/a relating to the translation remain the same.

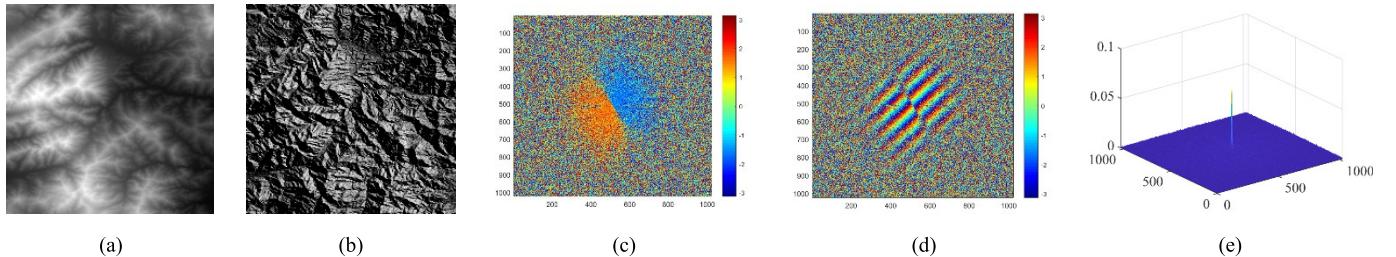


Fig. 2. Image pair and their PC result. (a) DEM image and (b) terrain shading optical image generated under the illumination of $\sigma = 10^\circ$ and $\tau = 60^\circ$. (c) and (d) PC cross power spectrum $Q(\omega, \theta)$ without and with image shift. (e) IFT of $Q(\omega, \theta)$.

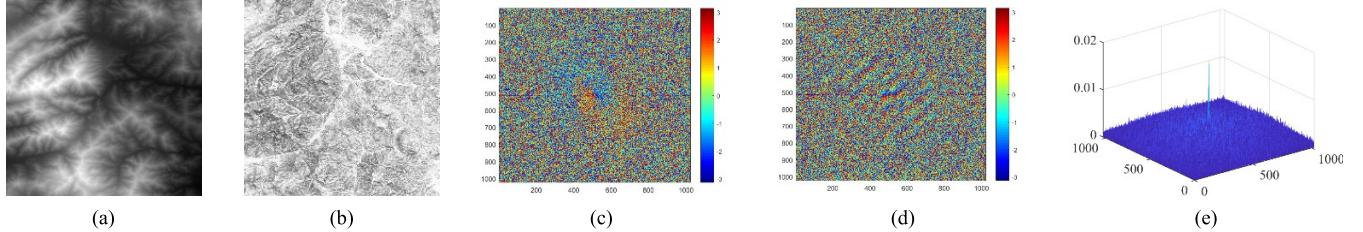


Fig. 3. Image pair and their PC result. (a) DEM image and (b) terrain shading optical image generated under the illumination of $\sigma = 5^\circ$ and $\tau = 60^\circ$. (c) and (d) PC cross power spectrum $Q(\omega, \theta)$ without and with image shift. (e) IFT of $Q(\omega, \theta)$.

As $\sin \sigma$, Z_d , and $\cos(\theta - \sigma)$ are all real numbers, the IFT of $Q(\omega, \theta)$ is still approximate to Dirac delta as shown in Fig. 2(e)

$$\text{IFT}(Q(\omega, \theta)) = \delta(x - a, y - b). \quad (9)$$

2) *Effect of Topography-Dependent Term*: When elevation angles σ are small, the value of $\cos \sigma$ is close to 0 and the term $\sin \sigma \approx 1$, so the PC cross power spectrum $Q(\omega, \theta)$, shown in Fig. 3(c), can be expressed as

$$Q(\omega, \theta) = \frac{F_H^*(\omega, \theta)}{Z_d}. \quad (10)$$

Here, $Q(\omega, \theta)$ becomes illumination independent, which means that the phase angle in the PC cross power spectrum is determined only by topography, irrelevant of illumination condition.

The cross power spectrum of $I(x, y)$ and $V_H(x, y)$ with image shift is

$$Q(\omega, \theta) = \frac{F_H^*(\omega, \theta)}{Z_d} e^{-i(au+bv)}. \quad (11)$$

Different from the effect of illumination term, the effect of topography term is not a systematically flipping of the fringe patterns, but a randomly flipping of the fringes as shown in Fig. 3(d). For a completely flat area, $V_H(x, y)$ is a constant and its FT $F_H^*(\omega, \theta)$ does not affect the PC Fourier spectrum $Q(\omega, \theta)$. For areas with topographic relief, the topographic vectors, $(\partial/\partial x)V_H(x, y)$ and $(\partial/\partial y)V_H(x, y)$, are variables depending on slope, and the results of $F_H^*(\omega, \theta)$ vary with topography. In other words, the impact of topography term degrades the power spectrum $Q(\omega, \theta)$, and however, the fringe density and orientation relating to the translation remain the same.

The IFT of $Q(\omega, \theta)$ is an approximation of the Dirac delta function, as shown in Fig. 3(e)

$$\text{IFT}(Q(\omega, \theta)) = \rho_H \delta(x - a, y - b) \quad (12)$$

where ρ_H is determined by the term of $\text{FT}((F_H^*(\omega, \theta))/Z_d)$.

3) *Effect of Terrain Albedo*: The above analysis is based on the Lambertian model that the albedo $r(x, y)$ of planetary terrain surface can be considered as a constant. This is true for most planetary exploration case, and however, there still contain situations that the terrain albedo cannot be considered as constant. This section will discuss the effect of terrain albedo on PC.

In the case terrain albedo values are variant, the PC cross power spectrum of optical image $F_I(\omega, \theta)$ and DEM $F_H(\omega, \theta)$ becomes

$$Q(\omega, \theta) = \left[\frac{\cos \sigma}{Z_d} \cos(\theta - \tau) + \frac{\sin \sigma}{Z_d} F_H^*(\omega, \theta) \right] R(\omega, \theta) \quad (13)$$

where $R(\omega, \theta)$ is the Fourier spectrum of albedo $r(x, y)$.

The cross power spectrum of $I(x, y)$ and $V_H(x, y)$ with albedo becomes

$$Q(\omega, \theta) = \frac{F_H^*(\omega, \theta) R(\omega, \theta)}{Z_d} e^{-i(au+bv)}. \quad (14)$$

Similar to the effect of topography-dependent term, the variation of albedo does not result in a systematical phase shift of the cross power spectrum, as it is mainly determined by $e^{-i(au+bv)}$. Thus, the effect of albedo variation degrades the power spectrum but will not change the density and orientation of the PC fringes.

In the general case, the negative correlation and randomly flipping introduced by topography-dependent term and illumination-dependent term will coexist in the PC cross power spectrum. However, none of the terms generate 2π wrapped

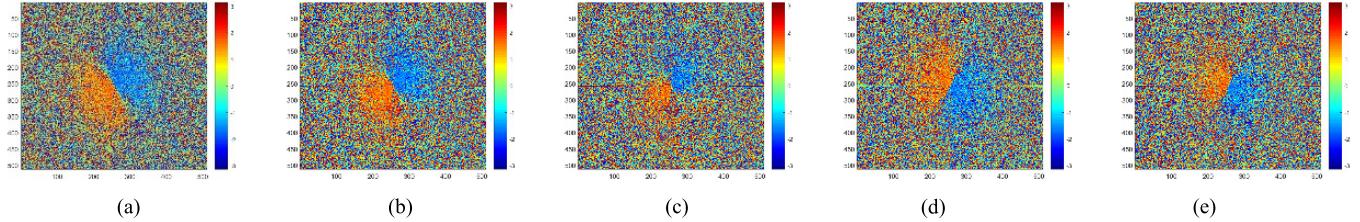


Fig. 4. Cross power spectrum maps and from the image pairs. This demonstrates different illumination conditions. (a) ($\tau = 60^\circ$, $\sigma = 10^\circ$). (b) ($\tau = 60^\circ$, $\sigma = 40^\circ$). (c) ($\tau = 60^\circ$, $\sigma = 70^\circ$). (d) ($\tau = 120^\circ$, $\sigma = 10^\circ$). (e) ($\tau = 120^\circ$, $\sigma = 45^\circ$).

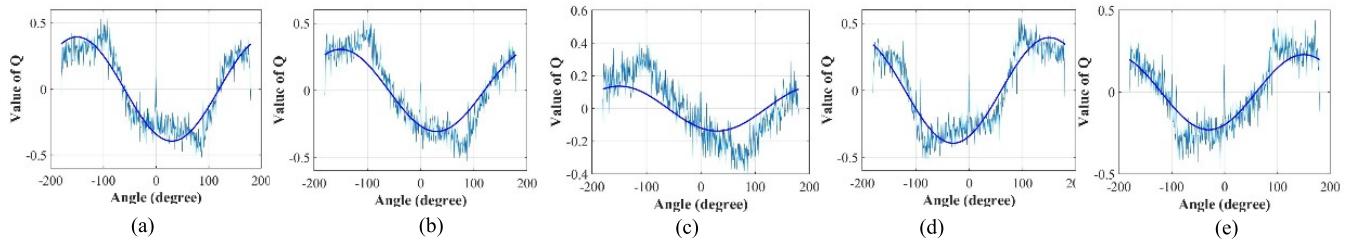


Fig. 5. Comparison between $Q(\omega, \theta)$ generated from image matching experiments (light blue) and $(\cos\sigma/Z_d)\cos(\theta - \tau)$ (dark blue). (a) ($\tau = 60^\circ$, $\sigma = 10^\circ$). (b) ($\tau = 60^\circ$, $\sigma = 40^\circ$). (c) ($\tau = 60^\circ$, $\sigma = 70^\circ$). (d) ($\tau = 120^\circ$, $\sigma = 10^\circ$). (e) ($\tau = 120^\circ$, $\sigma = 55^\circ$).

periodic fringe patterns as introduced by image translation $e^{-i(au+bu)}$. Therefore, although DEM data and optical images have a large difference in the spatial domain, their topography similarity can be identified in the frequency domain. This is the key property of PC-based image matching to make it robust to multimodal image registration.

Fig. 4 shows the PC cross power spectrums generated from the five image pairs, one optical image, and one DEM image, under different illumination conditions. It can be seen from the five figures that the whole cross power spectrums are divided into two angular parts: positive correlation parts shown by orange color and the negative correlation parts shown by blue color. This verified (6) that the value of $Q(\omega, \theta)$ is determined by the term $\cos(\theta - \tau)$. If the value of $\cos(\theta - \tau)$ is negative, the cross power spectrum will become negative. As shown in Fig. 4, under the same azimuth angle τ , the dividing line between positive correlation and negative correlation remains the same.

To further verify the relationship between $Q(\omega, \theta)$ and $\cos(\theta - \tau)$, the cross power spectrums $Q(\omega, \theta)$ are plotted with respect to azimuth angle τ ranges from $[-\pi, \pi]$, shown by the light blue in Fig. 5(a)-(e). Meanwhile, the theoretical value of $Q(\omega, \theta)$ according to (6) is also plotted to be compared with the experiment results.

It can be concluded from Fig. 5 that, generally, the $Q(\omega, \theta)$ values are in accordance with the illumination term $\cos(\theta - \tau)$, as the theoretical relationship between illumination condition and cross power spectrum derived in (6) can be proved. There appears random noise in the experimental results of $Q(\omega, \theta)$, shown in light blue, and this is because of the existence of topography terms. In fact, according to (6), the PC cross power spectrum is a sum of both illumination term and topography terms weighted by elevation angle σ , and thus, when elevation angles σ are small, the illumination term is dominant and

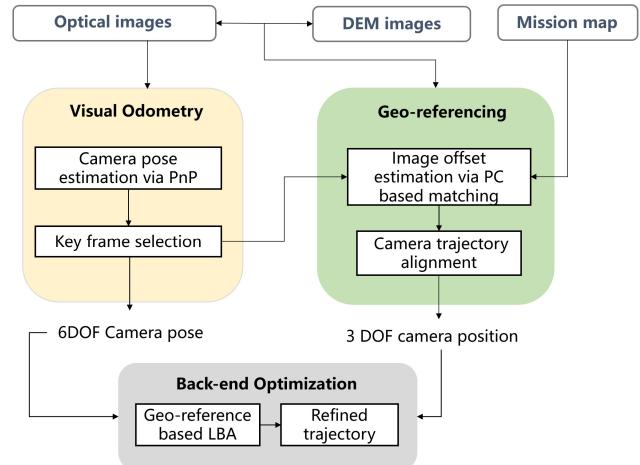


Fig. 6. Pipeline of the proposed terrain aided geo-referencing for planetary UAV optical navigation approach.

the values of $Q(\omega, \theta)$ are largely determined by the value of $\cos(\theta - \tau)$, as shown in Fig. 5(a) and (d). With the increase of elevation angle σ , the topography term become dominant, and thus, the effect of $\cos(\theta - \tau)$ becomes weak as shown in Fig. 5(c) that the blue curves become flat compared with other curves shown in Fig. 5.

IV. TERRAIN AIDED OPTICAL NAVIGATION

Based on the robustness of PC to multimodal image matching, a terrain aided optical navigation approach is proposed. The pipeline of the proposed optical navigation approach is shown in Fig. 6, which can be generally divided into three parts: **visual odometry**, **geo-referencing**, and **backend optimization**. Visual odometry estimated the **six-DOF camera**

pose via feature tracking. Then, keyframes are selected and localized to a reference map based on PC image matching. The trajectories estimated by visual odometry and geo-referencing are aligned and fused into a backend optimization. A new cost function is proposed for the LBA algorithm, which fuses the relative pose estimation from visual odometry and absolute pose estimation from geo-referencing to achieve the optimal pose estimation.

A. Visual Odometry Based on Frame-to-Frame Feature Tracking

Visual odometry is carried out via frame-to-frame feature tracking. The 2-D coordinate of feature points is determined via ORB features extraction and tracking [32]. Let $m_k^i = [u^i, v^i]$ to be the coordinates of the i th 2-D feature points on the image I_k at time k .

The initial 3-D point $M_j = [x_j, y_j, z_j]$ can be reconstructed using triangulation from the 2-D corresponding points via frame-to-frame matching.

The relationship between 2-D and 3-D points can be obtained via projection function $\pi(\cdot)$

$$m_k^i = \pi(P_k, M_j) \quad (15)$$

where $P_k \in \text{SE}(3)$ is the camera 3-D pose of image I_k

$$\text{SE}(3) = \left\{ T = \begin{bmatrix} R & t \\ 0^T & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4} \mid R \in \mathbb{R}^{3 \times 3}, t \in \mathbb{R}^{3 \times 1} \right\}. \quad (16)$$

The camera pose P_k can be estimated with a known 3-D map point M_j and 2-D feature point m_k^i via efficient PnP (EPnP).

Then, the camera trajectory t_{vo} can be recovered defined as a serial of camera coordinate with the increase of time k

$$t_{vo} = \{S_0^t, S_1^t, S_2^t, \dots, S_k^t\} \quad (17)$$

where $S_k^t = [X_k^t, Y_k^t, Z_k^t]$ is calculated from camera pose P_k as

$$S_k^t = S_{k-1}^t + t_k. \quad (18)$$

In visual odometry, the camera trajectory t_{vo} is in the local reference frame, which takes the first camera frame position S_0 as coordinate origin. The visual odometry does not solve the problem of estimation rover's position in a global map coordinate system, which can be solved via geo-referencing.

As it is not computationally efficient to take every image frame for geo-referencing, several keyframes I_k^i are selected for geo-referencing. The keyframe selection is based on the following criterion: 1) the time interval for two keyframes should not be too small, and otherwise, there is no point for the keyframe selection and 2) keyframes are opted to be salient with rich features, so as to be matched with a reference map.

Thus, in this article, we adopted the keyframe selection strategy as in ORB-SLAM2.

- 1) More than 20 frames have passed from the last keyframe.
- 2) There are more than 50 feature points in the current frame.
- 3) Current frame tracks less than 90% more points than the reference keyframe.

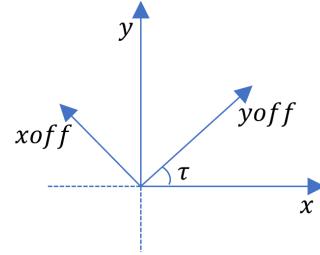


Fig. 7. Relationship between the current coordinate system and the map coordinate system.

B. PC for Terrain-Based Geo-referencing

The main purpose of geo-referencing is to provide the absolute position for image frames in a global map coordinate system. As proved in Section III, PC is able to correlate the topographic feature between multimodal images. Another merit of PC for absolute navigation is that PC is a noniterative matching algorithm. This means that the image offset can be directly calculated without roaming searching even if target UAV images and the reference terrain shading images only share a small part of the overlapping region. Thus, the PC-based image matching can largely shorten the matching time, which will ensure real-time optical navigation.

The offset between the current UAV image frame I_k and the reference terrain shading images R_k can be obtained via PC-based image matching

$$(X_{\text{off}}, Y_{\text{off}}) = \text{PC}(I_k, R_k). \quad (19)$$

Because the current coordinate system of the UAV image frame I_k depends on the current azimuth angle τ_k , it is essential to convert the current coordinate system into the global map coordinate system, as shown in Fig. 7.

The offset $(\Delta X_k, \Delta Y_k)$ between current UAV images and the reference terrain shading images can be calculated via the following equation:

$$\begin{bmatrix} \Delta X_k \\ \Delta Y_k \end{bmatrix} = \begin{bmatrix} -\sin \tau_k & \cos \tau_k \\ \cos \tau_k & \sin \tau_k \end{bmatrix} \begin{bmatrix} X_{\text{off}} \\ Y_{\text{off}} \end{bmatrix}. \quad (20)$$

The reference terrain shading images R_k are generated using (4) under presumably illumination conditions of given solar radiation intensity L in the direction of azimuth angle τ and elevation angle σ . To speed up the image registration, the DEM data are cropped to local DEM according to the planned flight path T_{pl}

$$T_{\text{pl}} = \{S_0^R, S_1^R, S_2^R, \dots, S_k^R | S_k^R = [X_k^R, Y_k^R, Z_k^R]\}. \quad (21)$$

According to the offset $(\Delta X_k, \Delta Y_k)$ between current UAV images and the reference terrain shading images, the camera trajectory in the global map coordinate system $T_{vo} = \{S_0^T, S_1^T, S_2^T, \dots, S_k^T | S_k^T = [X_k^T, Y_k^T, Z_k^T]\}$ can be estimated via the following equation:

$$\begin{aligned} X_k^T &= X_k^R + \text{GSD} \cdot \Delta X_k \\ Y_k^T &= Y_k^R + \text{GSD} \cdot \Delta Y_k \\ Z_k^T &= \text{DEM}(X_k^R, Y_k^R) + H \end{aligned} \quad (22)$$

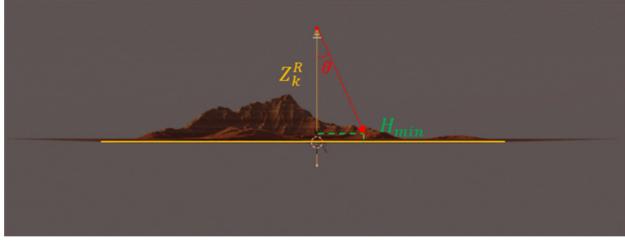


Fig. 8. Illustration of GSD with respect to flight height and topography.

where $\text{DEM}(X_k^R, Y_k^R)$ is the elevation of the ground surface at the X_k^R - and Y_k^R -coordinates and H is the height of the UAV to the ground surface. GSD is the ground sample distance. Note that the UAV flies relatively low at the Martian surface, the GSD will change with the change of flight height and surface topography as shown in Fig. 8

$$\text{GSD} = 2 \frac{Z_k^R - H_{\min}}{s} \tan \theta \quad (23)$$

where H_{\min} is the lowest altitude position in the field of view of the camera at this altitude, which is related to the focal length of the camera, and s is the size of the UAV image defined as

$$s = \max(w, h) \quad (24)$$

where w and h are the width and height of the UAV images, respectively. For example, if the image size is 720×480 pixels, the value of s is 720.

From geo-referencing based on PC matching, the camera trajectory T_{vo} of the keyframe is recovered in the global map coordinate system, while in visual odometry, the camera trajectory t_{vo} is estimated in the local reference coordinate system. The two trajectories are required to be transformed into the same coordinate system in order to be fused into an LBA. The transformation between the two trajectories can be calculated as

$$T_{vo} = H_{Tt} t_{vo} \quad (25)$$

where $H_{Tt} \in \text{SE}(3)$.

C. Keyframe Pose Refinement Based on LBA

The results of visual odometry and geo-referencing are complementary. The error accumulation or “trajectory drift” can be corrected by the result of geo-referencing via global matching, while visual odometry can produce six-DOF camera poses compared to three-DOF camera poses estimated from geo-referencing. Thus, LBA is performed to fuse the visual odometry and geo-referencing results into more reliable and accurate camera poses.

Bundle adjustment (BA) can be defined as the problem of simultaneously refining the 3-D coordinates describing the scene geometry, the parameters of the relative motion, and the optical characteristics of the camera(s) employed to acquire the images, according to an optimality criterion involving the corresponding image projections of all points [33]. Instead of taking all 3-D points and camera poses for optimization,

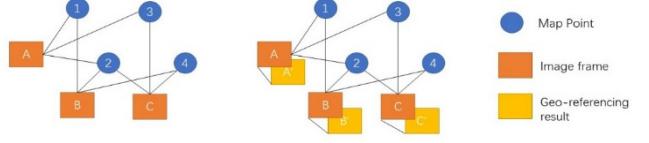


Fig. 9. Graph network of traditional BA and the proposed geo-referencing-based LBA.

which is time-consuming, LBA takes several adjacent camera positions, such as 15 frames, for optimization. This will largely increase the optimization speed during real-time navigation.

LBA can be done for a fixed interval, such as every 15 or 20 image frames. The problem is that the error propagation in SLAM is not a linear increase with the distance. For example, the error will propagate faster when the UAV is turning or changing direction than keeping straightforward. There might also appear large localization errors when few features can be extracted from optical images.

Traditionally, LBA is achieved by minimizing the reprojection error between the image locations of observed and predicted image points. Fig. 9(a) shows the network structure of four 3-D points 1–4 and three images A–C; 3-D point 1 is seen in, point 2 is seen in A, B, and C, point 3 is seen in A and C, and point 4 is seen in B and C. Then, the cost function $f(x)$ is defined as

$$f(x) = \frac{1}{2} \sum_{i=0}^M \sum_{j=0}^N e_{i,j}^T w_i e_{i,j} \quad (26)$$

where $e_{i,j} = \pi(P_k, M_j) - m_k^i$ is the reprojection error of 3-D point M_j on the image I_i and w_i is the confidence of observation m_k^i .

In this article, the locations of an image can be estimated via geo-referencing, as shown in Fig. 9, and thus, a new cost function $f(x)$ is proposed to fuse the geo-referencing result into LBA

$$f(x) = \frac{1}{2} \sum_{i=0}^M \sum_{j=0}^N \left(e_{i,j}^T w_i e_{i,j} + e_i^G w_i^G e_i^G \right) \quad (27)$$

where w_i^G is the confidence of geo-referencing and $e_i^G = [e_x, e_y, e_z]$ are deviations between geo-referencing and visual odometry in the x -, y -, and z -directions that are defined as

$$\begin{aligned} e_x &= |X_k^T - X_k^{T'}| \\ e_y &= |Y_k^T - Y_k^{T'}| \\ e_z &= |Z_k^T - Z_k^{T'}| \end{aligned} \quad (28)$$

where (X_k^T, Y_k^T, Z_k^T) is the calculated camera position from geo-referencing and $(X_k^{T'}, Y_k^{T'}, Z_k^{T'})$ is the camera position from visual odometry.

As the cost function $f(x)$ is nonlinear, the Taylor series expansion is applied to $f(x)$ for linearization. Then, the optimization is achieved by Levenberg–Marquardt (LM).

V. EXPERIMENT RESULTS

A. Mars UAV Image Dataset

In this section, nine simulated Mars scenes, as shown in Fig. 10, were generated via Blender, a free, and open-source

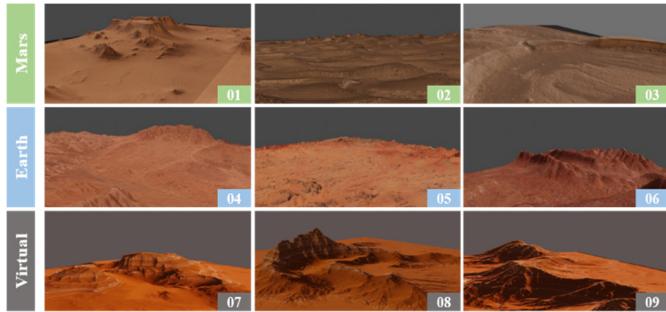


Fig. 10. Nine simulated Mars scenarios generated from different sources.

TABLE II
OPTIC STATISTICS SUMMARY OF NINE MARS SCENES

Scene	Terrain type	Data Sources	DEM Resolution (m/pixel)	Scene coverage (km)	Resize factor
01	highlands	HiRise	1.00	5.9×8.5×0.5	10
02	dunes	HiRise	1.00	5.8×9.6×0.2	10
03	crater	HiRise	1.00	6.2×9.1×0.3	10
04	highlands	Google Map	8.56	10.3×10.3×1.5	20
05	gravel	Google Map	16.30	19.6×19.6×3.9	38
06	mountain	Google Map	7.64	9.1×9.1×0.5	18
07	canyon	simulation	0.07	0.6×0.6×0.07	1
08	valley	simulation	0.07	0.6×0.6×0.08	1
09	mountain	simulation	0.15	1.2×0.6×0.08	1

3-D scene simulation and rendering engine.² The data source of the DEM model and image texture includes Earth observation data, Mars orbital imagery, and a simulated Mars-like 3-D terrain model. Then, based on the nine scenes, 23 UAV image sequences were generated from different trajectories. The UAV image size is 480 × 480 pixels, the focal length of the camera and the sensor size are both set to 36 mm, the elevation and azimuth angles of sunlight in the environment are (60°, 0°), and the light strength is 10. The Mars UAV dataset contains 40 200 images in total.

The simulated Mars scenes are shown in Fig. 10. 01–03 scenes are generated from the High Resolution Imaging Science Experiment³ (HiRise) equipped on the Mars Reconnaissance Orbiter with a ground resolution of 25 cm. The three scenes include the typical Martian topography such as slope streaks, dunes, and craters in Arabia Terra, Nili Patera, and Breachin. 04–06 scenes are Mars-like terrain models generated from Google Map, which are scenes from Actama Desert, Jordan Desert, and Chaidan Cliff. 07–09 scenes are three simulated 3-D terrain models of canyons, valleys, and mountains.

Statistics details of the nine Mars scenes are summarized in Table II. As there are three types of data sources for Mars 3-D scene simulation, our dataset represents different reflect Mars scenes at different scales. The terrain model from GoogleMap covers a large area of terrain models with rich types of topographic features, while the simulated terrain models (07–09 scenes) contain detailed topographic features. The comparison of different scene scales can be seen in

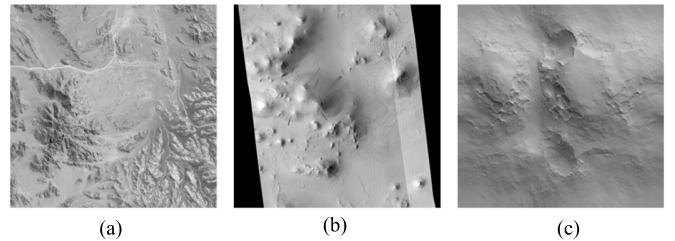


Fig. 11. Different scales in Mars scenes. (a) Google Earth. (b) HiRise. (c) Simulation.

Fig. 11. Using three Mars-like terrain models on the Earth also makes it possible for field trips in the future. Real UAV data can be taken on the Earth for accuracy assessment and algorithm performance enhancement.

To keep the size of nine scenes consistent for UAV data generation, all the DEM models are resized to about 600 × 600 m², together with texture image, and thus, the resize factor in Table II varies from one scene to another. For scenes 01–03, terrain texture is extracted from HiRise satellite imagery after orthorectification. For scenes 04–06, terrain texture is the surface texture extracted from GoogleMap. To simulate the color of Mars terrain, the colors of texture for 07–09 are manually changed to dark red.

The process of Mars 3-D terrain scene modeling and UAV image generation using Blender is shown in Fig. 12. First, we obtain the DEM and surface texture images of a scene through the three data sources: HiRise, Google Map, and simulation. Next, the DEM image and the surface texture image were imported into Blender, then create a terrain plane in Blender, and align the plane with the DEM and texture image. Moreover, since the value of each pixel on the DEM image represents the terrain height at that location, terrain elevation can be generated by Blender and then overlay the aligned surface texture image on the surface, and a preliminary scene model with surface textures is obtained. Finally, we set the properties of the terrain model such as reflectivity, transmittance, and color, as well as the altitude, azimuth, and intensity of the sun, and then, we set the focal length, resolution of the camera in Blender.

After establishing the 3-D scene model of Mars, the next step is to simulate the trajectories of UAVs and generate image sequences. Three UAV trajectories, circular, forward, and scanning, are generated and shown in Fig. 13. The circular trajectory simulates a touch-down looking on a specific region, the forward trajectory simulates a UAV moving forward directly, and the scanning trajectory simulates the quick mapping of a large Mars region. The whole flight time for the first trajectories is set as 900 s, and for the other two trajectories, it is set as 600 s. The UAV viewing angle is set as nadir view to ensure high mapping and localization precision. The camera frequency is set as 1 Hz. Finally, we can use Blender's rendering capabilities to obtain a sequence of images from the UAV's perspective.

B. Image Matching Between UAV Imagery and DEM

In this experiment, a UAV image sequence is used as target images, while another image sequence of DEM, generated

²<https://www.blender.org/>

³<https://hirise.lpl.arizona.edu/>

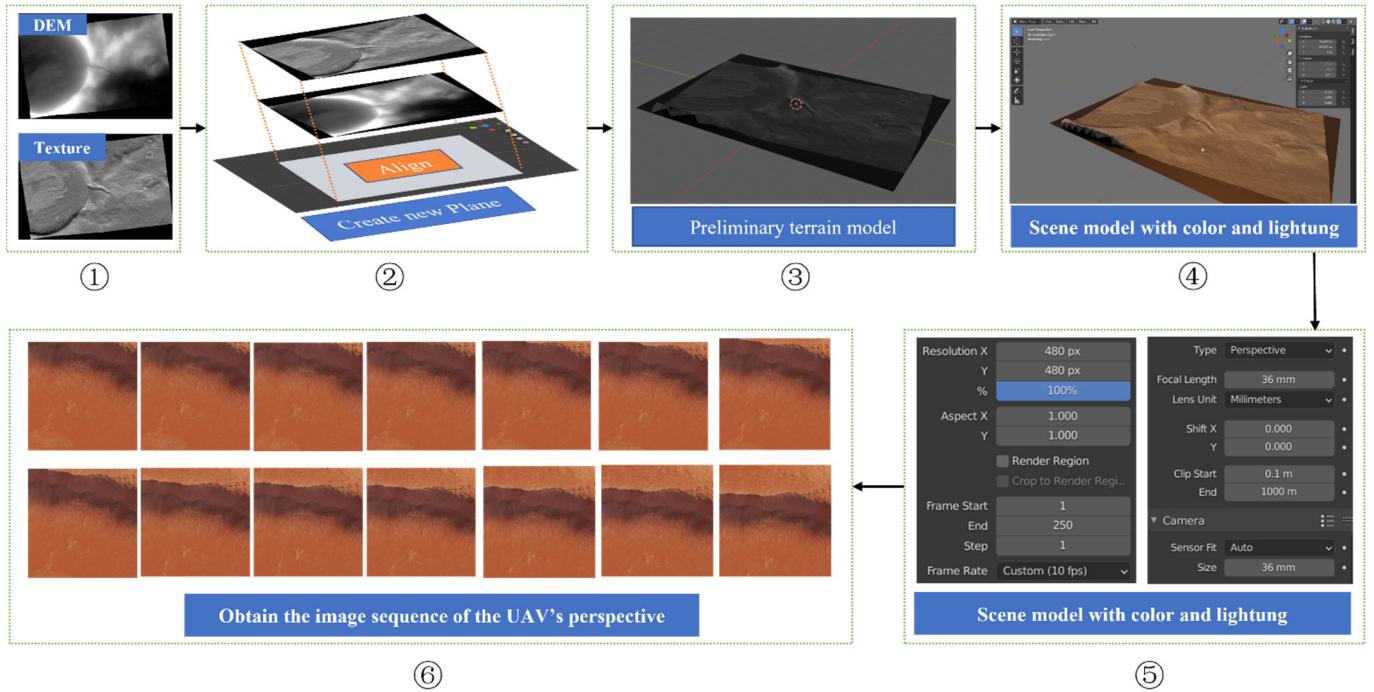


Fig. 12. 3-D Mars scene modeling and UAV image generation.

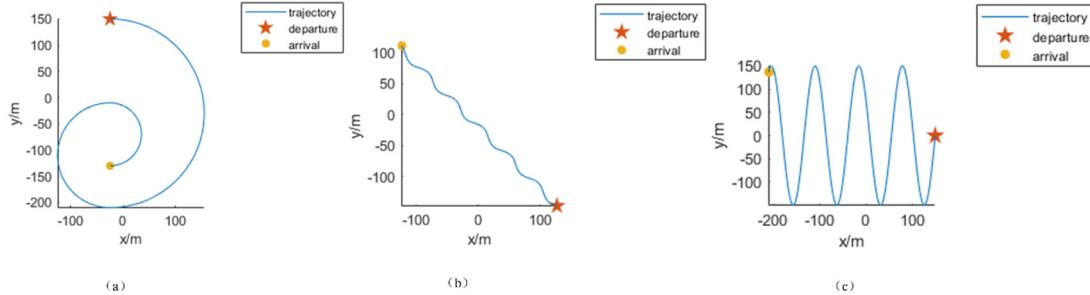


Fig. 13. Three flight trajectories of the drone: (a) circular trajectory, (b) forward trajectory, and (c) scanning trajectory (start–end).

TABLE III
IMAGE MATCHING ACCURACIES (PIXEL) USING SIFT
AND OUR METHOD (PC)

Algorithms	01	02	03	04	05	06	07	08	09
SIFT	--	0.31	--	--	0.28	--	4.39	0.29	--
Our method	0.22	0.07	2.09	0.05	0.06	0.21	3.12	0.03	0.13

using the same trajectory, is used as a reference. Since the same trajectory is used to generate the DEM and the optical image sequence, the ideal image offset is 0. Thus, the image offset between the optical images and the DEM terrain shading images can be used as a metric to evaluate the image matching accuracy.

For comparison, the image matching accuracies of a feature-based method, SIFT, are calculated. Table III lists the matching accuracies using PC and SIFT in the nine Mars scenes.

It can be seen in Table III that SIFT failed in most cases and only succeeded in four scenes. The feature correspondence

result of the grayscale images of the DEM image and the optical image in the nine scenes using SIFT is shown in Fig. 14. It can be seen from Fig. 14 that in most scenes, SIFT can hardly generate enough feature pairs due to the low texture in the Mars scene and large appearance variation between optical images and DEM terrain shading images.

The image matching performance of PC is superior to SIFT in all tested scenarios. As shown in Table III, except for scenes 3 and 7, the matching error of PC in all scenes is less than 1 pixel, and the average matching accuracy is 0.66 subpixel.

To further analyze the matching performance of PC in nine Mars scenes, error distributions of image matching are plotted in Fig. 15. It can be seen that the error distributions vary from scene to scene. For most scenes, the image offsets are distributed around 0, which can be modeled by the Gaussian distribution. However, for some cases, such as scenes 03 and 06, the error distribution can be hardly modeled by the Gaussian distribution. Fig. 15 also shows that although,

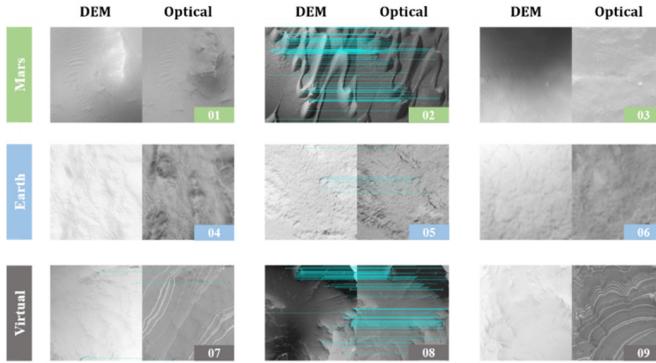


Fig. 14. SIFT-based image matching between optical UAV image sequence and the corresponding DEM terrain shaded images.

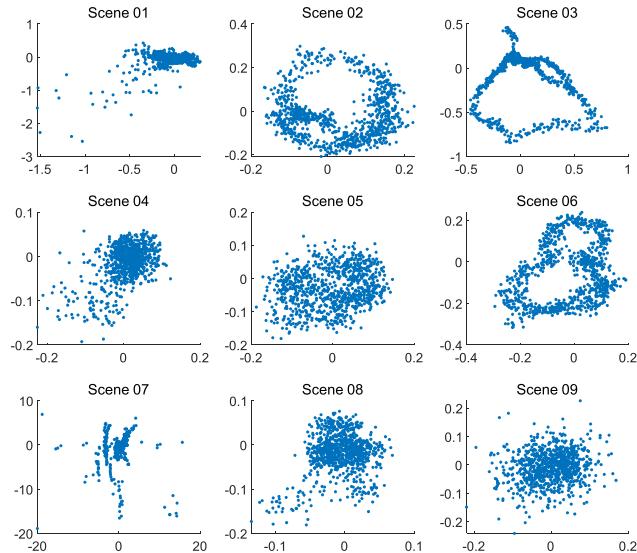


Fig. 15. Error distribution of image matching accuracy using PC in nine Mars scenes.

generally, PC achieves robust matching performance, it did fail in several cases, for example, in scene 07, for some images, the matching error can be as large as 20 pixels. This experiment paves the way for the PC-based optical navigation approach, which fuses geo-referencing with visual odometry for accurate and robust localization results.

C. Impact of Flight Translation and Rotation on Optical Navigation

In the actual flight of UAV, yaw is inevitable, which may cause the translation on the horizontal plane and the rotation of flight direction. In order to verify the robustness to flight translation and rotation, the following experiments, whose UAV image sequences and DEM image sequences are generated from the 3-D model of scene 8, are arranged in the circumstance of only translation, only rotation, and both translation and rotation.

1) *Impact of Flight Translation*: This experiment will investigate the impact of flight translation on geo-referencing results. According to the 3-D model of scene 8, Gaussian

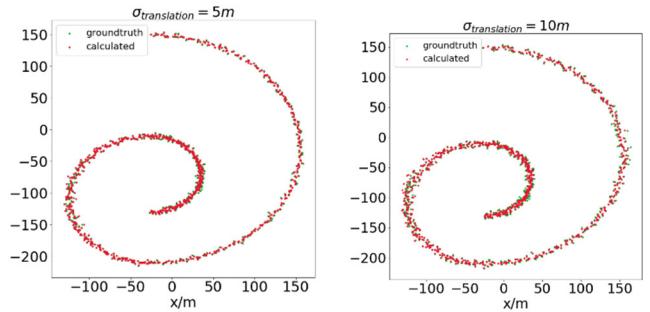


Fig. 16. Actual flight trajectories and calculated flight trajectories under different variances.

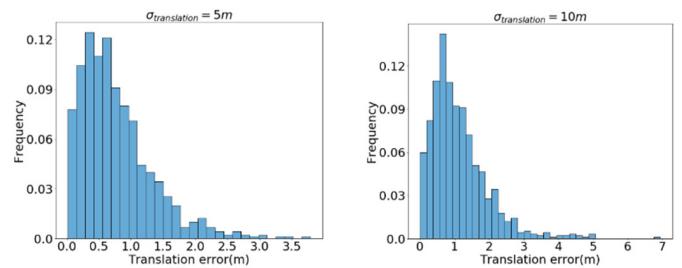


Fig. 17. Distribution of translation errors under different variances.

noises are added to the X - and Y -coordinates of DEM image sequences, with the variance of 5 and 10 m, respectively, and the mean of 0 m. Two groups of UAV image sequences are generated under these two different Gaussian noises, so as to simulate the phenomenon of UAV flight translation.

To assess the effect of the experiment, the coordinate values of the UAV image sequence are used as ground truth. According to (19)–(23), the pixel deviation of PC can be converted to the coordinate deviation relative to the DEM image sequence to get the coordinates calculated value. The Euclidean distance between coordinates calculated value and ground truth of Euclidean distance can be used as the error to measure the experimental results.

The comparison between the calculated trajectory and the actual trajectory is shown in Fig. 16. It can be seen that the calculated trajectory has a high similarity with the actual trajectory. The error distribution of each picture is shown in Fig. 17. When the variance is 5 m, the average error is 0.77 m; when the variance is 10 m, the average error is 1.13 m.

The experiment shows that under the influence of different Gaussian noises, the accuracy of geo-referencing is acceptable. Even if the actual trajectory does not match the preset one, it can restore the real position and has the **robustness to flight translation**.

2) *Impact of Flight Rotation*: This experiment will investigate the impact of flight rotation on georeference results. According to the 3-D model of scene 8, Gaussian noises are added to the yaw angle of DEM image sequences, with the variance of 5° and 10° and the mean of 0° . Two groups of UAV image sequences are generated under these two different Gaussian noises to simulate the phenomenon of UAV flight translation.

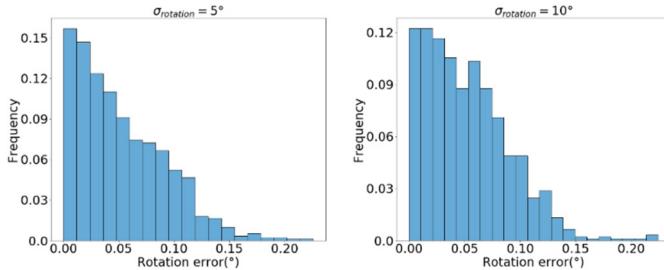


Fig. 18. Distribution of rotation errors under different variances.

To assess the effect of the experiment, the yaw angle used to generate the UAV image sequence is used as ground truth. According to (3), the rotation deviation of the yaw angle relative to the DEM image sequence can be calculated to obtain the calculated value of the yaw angle, and the absolute value of the difference between the yaw angle calculation value and the ground truth can be used as the error to measure the experimental effects.

The error distribution of each picture is shown in Fig. 18. When the variance is 5° and 10° , the average error is 0.05° . It is probably because the georeference has a fair performance under small rotation.

The experiment shows that under the influence of different Gaussian noises, the accuracy of geo-referencing is acceptable. It can calculate the real value of the yaw angle and it has the robustness to flight rotation.

3) Impact of Flight Translation and Rotation: This experiment will investigate the combined impact of flight translation and rotation on georeference results. According to the 3-D model of scene 8, Gaussian noises are added to the X - and Y -coordinates of DEM image sequences, with the variances of 5 and 10 m, respectively, and the mean of 0 m. Gaussian noises are also added to the yaw angle of DEM image sequences, with the variance of 5° and 10° and the mean of 0° . Two groups of UAV image sequences are generated under two different kinds of Gaussian noises. In this way, the flight translation and rotation caused by various factors in the real flight of UAV are simulated.

To assess the effect of the experiment, the yaw angle used to generate the UAV image sequence and the coordinate values of the UAV image sequence are used as ground truth. The calculation process is divided into two steps. According to (3), the rotation deviation of the yaw angle relative to the DEM image sequence is calculated. Based on the rotation deviation, the calculated value of the yaw angle is obtained. Then, the pixel deviation obtained by phase matching can be converted into the coordinate deviation relative to the DEM image sequence by referring to formulas 16–20 to obtain the calculated value of coordinates. The absolute value of the difference between the calculated yaw angle and the ground truth is the error of the yaw angle, and the Euclidean distance between the calculated coordinates and the ground truth is the error of the coordinates.

The comparison between the calculated and the actual trajectory is shown in Fig. 19. It can be seen that the

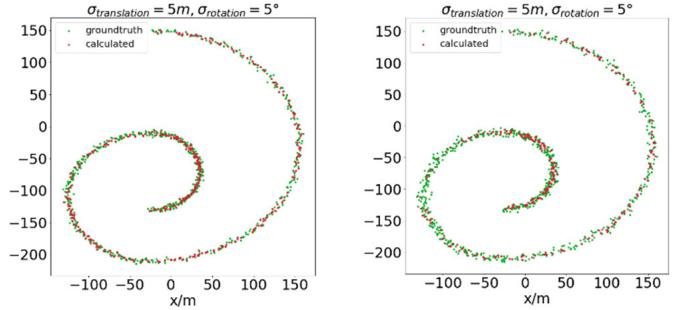


Fig. 19. Actual flight trajectories and calculated flight trajectories under different variances.

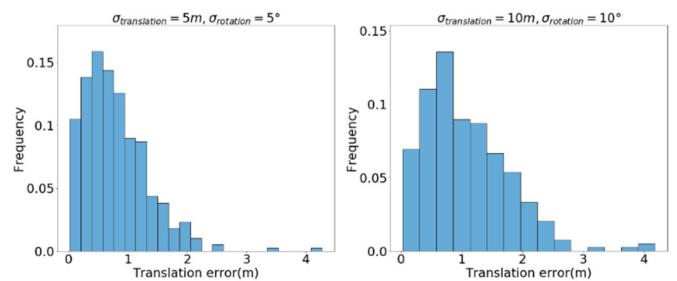


Fig. 20. Distribution of translation errors under different variances.

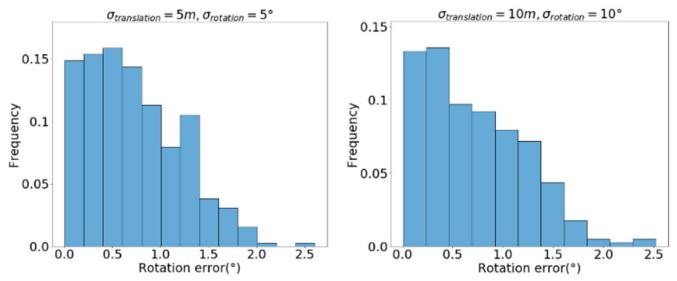


Fig. 21. Distribution of rotation errors under different variances.

similarity between the calculated and the actual trajectory is high, but there is incompleteness in the calculated track. Georeferencing cannot guarantee the reliability of each image in the sequence under complex conditions. By filtering the peak value of phase matching and changing the size of Hamming window, more reliable picture frames can be retained and the average error can be reduced.

The error distribution of each picture is shown in Figs. 20 and 21. When the variance of Gaussian noise is 5 m and 5° , the average translation error is 0.79 m and the average rotation error is 0.73° ; when the variance of Gaussian noise is 10 m and 10° , the average translation error is 1.16 m and the average rotation error is 1.04° .

After a series of experiments, as shown in Table IV, the accuracy of georeference is affected to some extent by adding Gaussian noise to coordinate and yaw angle. Compared with only yaw angle noise, the yaw angle error changes more with the influence of translation. It may be because, under the joint action of rotation and translation, the common view area of

TABLE IV

AVERAGE TRANSLATION ERROR AND AVERAGE ROTATION ERROR UNDER THE INFLUENCE OF DIFFERENT GAUSSIAN NOISES

Average error	$\sigma_t = 5\text{m}$	$\sigma_t = 10\text{m}$	$\sigma_r = 5^\circ$	$\sigma_r = 10^\circ$	$\sigma_t = 5\text{m}$	$\sigma_t = 10\text{m}$
	$\sigma_r = 5^\circ$	$\sigma_r = 10^\circ$				
Translation(m)	0.77	1.13	--	--	0.79	1.16
Rotation($^\circ$)	--	--	0.05	0.05	0.73	1.04

TABLE V

TRAJECTORY ERROR UNDER FIVE DIFFERENT CONFIDENCES

confidence	max/m	median/m	min/m	mean/m	rmse/m	std/m
0	6.10	1.76	0.28	1.99	2.22	0.99
1	4.58	0.88	0.08	0.94	1.07	0.52
10	2.86	0.30	0.007	0.29	0.38	0.24
100	2.12	0.32	0.05	0.41	0.51	0.31
1000	2.80	0.31	0.02	0.39	0.48	0.28

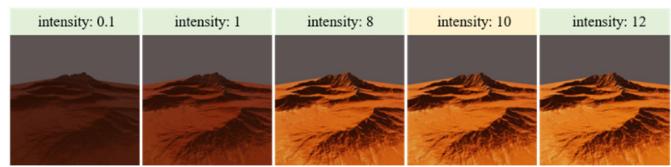


Fig. 22. Appearance of scene 8 under different illumination intensities.

the two images decreases, and each image has many new textures that the other one does not have, which will greatly affect the calculation results of PC, resulting in very large calculation results. Even if the true value of yaw angle cannot be completely restored, the mean error of 1° is still very small for the calculation of translation. Besides, geo-referencing only needs to provide fine translation results and not all translation and rotation results. It is reasonable to take geo-referencing as the terrain constraint of the visual odometer because of the robustness of the translation and rotation of the flight based on PC.

D. Impact of Illumination Condition on Optical Navigation

Lighting is one of the major environmental changes on Mars, which significantly alters the quality of UAV images and thus influences navigation accuracy. The following experiments begin to explore the robustness of PC to lighting conditions in the aspects of intensity, azimuth angle, and elevation angle.

1) *Change of Illumination Intensity*: This experiment will investigate the effect of lighting strength on geo-referencing results. Five UAV image sequences are generated from the 3-D model of scene 8 under five different illumination intensities from 0.1 to 12, as shown in Fig. 22.

The image sequence with the illumination intensity of 10 is set as the presumable illumination condition for terrain shading image generation, and other parameters are the same among different image sequences except for the illumination intensity. Then, we test the matching accuracy of PC with respect to illumination intensity using the same metric as in Section V-B.

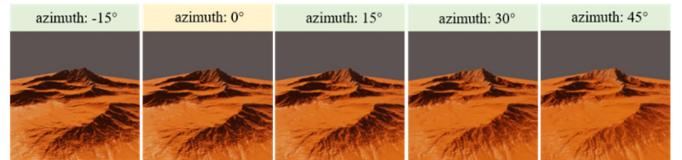


Fig. 23. Appearance of scene 8 under different illumination azimuths.

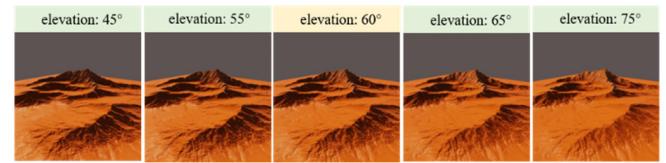


Fig. 24. Appearance of scene 8 under different illumination elevations.

The image matching accuracies of PC under different illumination intensities are shown in Fig. 25. It can be seen from Fig. 25 that when the illumination intensity is 0.1, the maximum matching error is 0.129 subpixel, and with the lighting strength approaching 10, the error reduces to 0.047 subpixel. Fig. 26 shows the geo-referencing error using PC and has a similar trend. The maximum localization error dropped from 0.032 to 0.015 m with the lighting strength approaching 10.

This experiment demonstrates that the accuracy of the geo-referencing is higher when the difference in light intensity between the presumably DEM image and the real optical image is smaller. However, even if the presumably light intensity is not accurate, PC can still produce an accurate geo-referencing result.

2) *Change of Illumination Azimuth*: This experiment will investigate the effect of illumination azimuth on geo-referencing results. Five UAV image sequences are generated from the 3-D model of scene 8 under five different illumination azimuths from -15° to 45° , as shown in Fig. 23. The presumable azimuth angle is 0° .

The accuracies of image matching and geo-referencing using PC are shown in Figs. 27 and 28. Similar to the experiment result in Section V-C1 that the matching error reduces to 0.047 subpixel when the light azimuth angle is precisely estimated.

3) *Change of Illumination Elevation*: This experiment will investigate the effect of illumination elevation on geo-referencing results. Five UAV image sequences are generated from the 3-D model of scene 8 under five different illumination elevation angles from 45° to 75° , as shown in Fig. 24. The presumable elevation angle is 60° .

The accuracies of image matching and geo-referencing using PC under different illumination elevations are shown in Figs. 29 and 30. The average image matching accuracy is 0.062 subpixel, and the average geo-referencing accuracy is 0.089 m.

The lighting condition may not always be accurately estimated for every scheduled task. However, the experimental results prove that PC-based geo-referencing has robustness to the uncertainty of lighting conditions. Therefore, it is

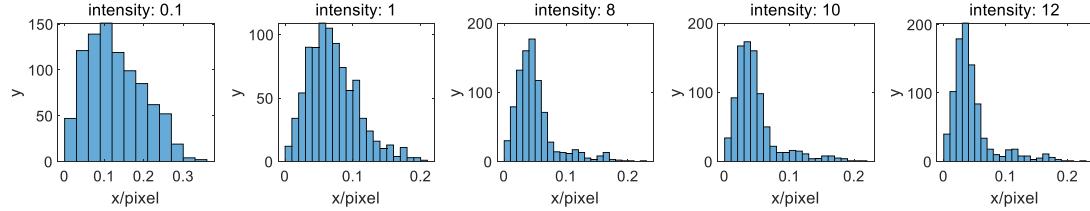


Fig. 25. Pixel offsets of optical images under different illumination intensities.

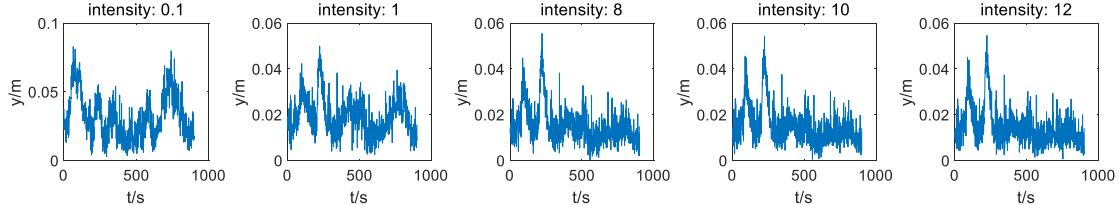


Fig. 26. Absolute errors of geo-referencing under different illumination intensities.

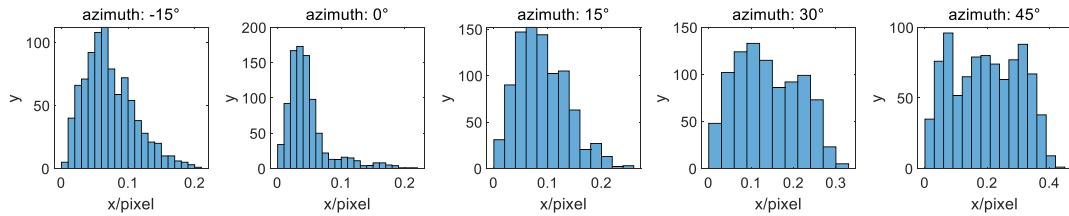


Fig. 27. Pixel offsets of optical images under different illumination azimuths.

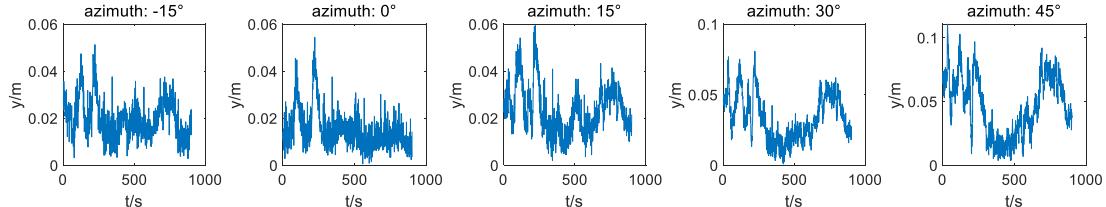


Fig. 28. Absolute errors of geo-referencing under different illumination azimuths.

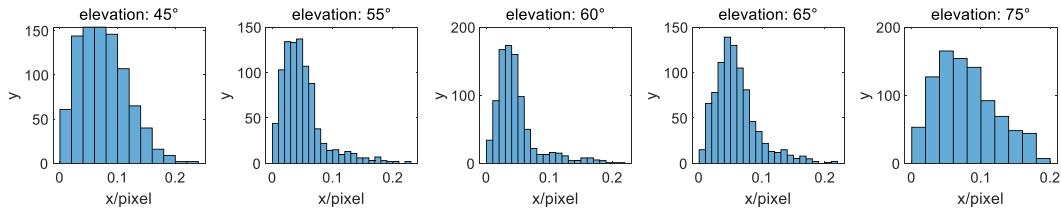


Fig. 29. Pixel offsets of optical images under different illumination elevations.

quite feasible to use the geo-referencing results as the terrain constraint in the visual odometer.

E. Terrain Aided Visual Odometry Experiment

1) *Influence of the Confidence of Terrain Constraints on the Results of LBA:* The above experiment results confirm the robustness and accuracy of PC-based geo-referencing. In this experiment, we will further investigate the effect of the terrain constraints under different confidence values.

The terrain aided VO experiments are carried out using the UAV data of scene 8. The optical images have been generated when the light intensity is 12, the azimuth angle is 0°, and the altitude angle is 60°. A total of six groups of experiments are set up with the confidence level of 0, 1, 10, 100, and 1000. The experiment with a confidence of 0 is equivalent to terrain constraint-free. The EVO tool is used to evaluate the accuracy of calculated trajectories.⁴ The absolute

⁴<https://github.com/MichaelGrupp/evo>

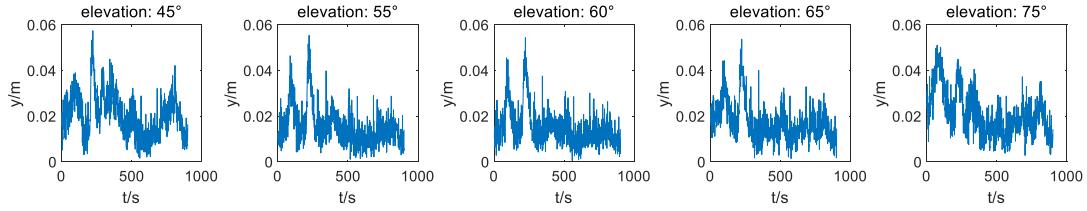


Fig. 30. Absolute errors of geo-referencing under different illumination elevations.

TABLE VI
POSITIONING ACCURACY UNDER THREE FLIGHT TRAJECTORIES

Trajectory	confidence	max/m	min/m	mean/m	rmse/m	std/m
a	ORB-SLAM2	5.61	0.35	1.92	2.13	0.93
	Our Method	2.86	0.007	0.29	0.38	0.24
b	ORB-SLAM2	1.93	0.04	0.55	0.63	0.30
	Our Method	1.29	0.03	0.24	0.28	0.15
c	ORB-SLAM2	3.15	0.16	1.15	1.22	0.41
	Our Method	2.24	0.01	0.36	0.48	0.32

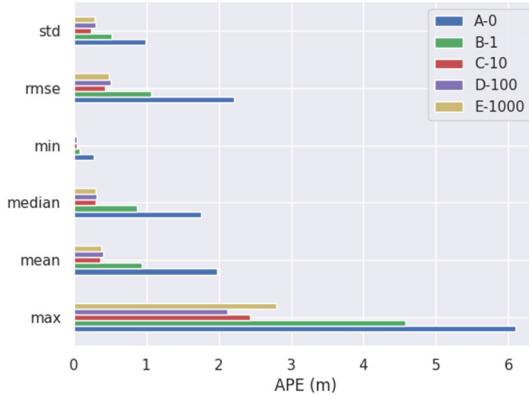


Fig. 31. Absolute pose error (APE) with respect to five different confidences.

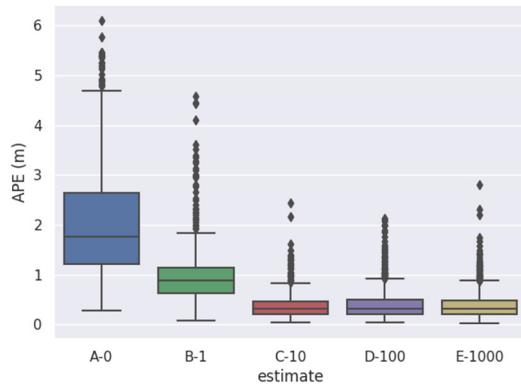


Fig. 32. Box plot of APE with respect to five different confidences.

errors of the trajectories under five different confidences are shown in Table V, Figs. 31 and 32.

It can be seen from Fig. 32 that when the confidence is equal to 10, the average error, root-mean-square error, and standard

deviation are all the smallest. In addition, it can be seen that the accuracy increases significantly when the confidence level is from 0 to 10. This shows that the terrain constraint based on PC has played a significant role to increase the whole navigation accuracy. However, the increased accuracy is not very obvious when the confidence is from 10 to 1000, indicating that the optimal confidence value is around 10.

2) *Comparison Results of Different Flight Trajectories*: This experiment will analyze the effect of UAV trajectories on the accuracies of navigation. Fig. 33 shows the navigation errors of the estimated trajectories and the ground-truth trajectories. The color bar on the right represents the absolute positioning error of the corresponding color. The statistic errors can be seen in Table VI.

As shown in Fig. 33, most points in calculated trajectories are displayed in blue with a low error level. Table VI shows that compared to ORB-SLAM2 which produces 1.3-m average localization accuracy, the proposed method is able to achieve the average localization accuracy of 0.3 m. ORB-SLAM2 has the lowest localization accuracy in the circular path, which may owe to the loop-closure error in circle-like trajectories, the algorithm may mistakenly consider the current position returns to the starting point, while the actual trajectory is not. It is worthwhile to mention that the proposed method is not sensitive to trajectories, as the optimization is not based on loop closure, but based on geo-referencing.

3) *Terrain Aided VO Experiments in All Scenarios*: This experiment will investigate the performance of the proposed method for all nine scenes with different terrain types and scene scales. The lighting conditions vary from the scene and can be seen in Table VII. As proved in Section V-C2, the proposed algorithm is not sensitive to flight trajectories, and thus, we use circular trajectories in all nine scenes.

We compare the proposed method with four state-of-the-art SLAM algorithms, ORB-SLAM2, ORB-SLAM3, SURF-SLAM, and SuperPoint-SLAM. The localization accuracies,

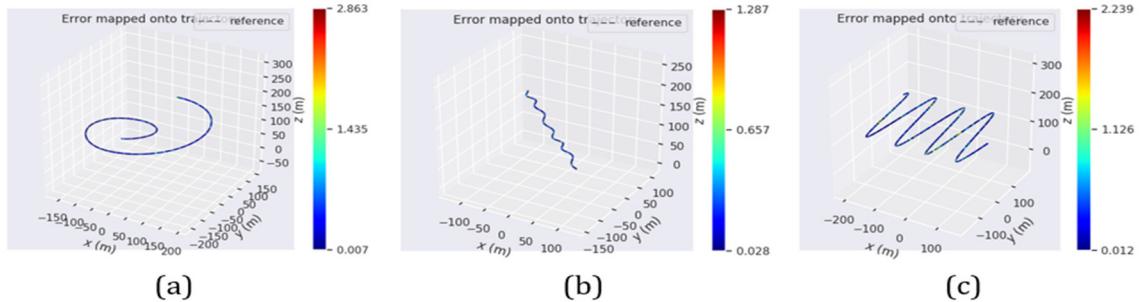


Fig. 33. Navigation error of estimated trajectory and true trajectory in three cases. (a) Circular trajectory, (b) Forward trajectory, and (c) Scanning trajectory.

TABLE VII
LIGHTING CONDITIONS OF THE NINE SCENES IN THE EXPERIMENT

Scene	Light strength	Light azimuth/°	Light elevation/°
01	5	30	60
02	2	0	60
03	5	0	60
04	5	30	75
05	2	0	30
06	5	0	30
07	5	0	90
08	10	0	60
09	10	0	60

including maximum and mean errors, average errors, and standard deviation, are shown in Table VIII. The bold numbers in Table VIII stand for the best performance among the relevant metrics, while dashes mean that the SLAM method fails to complete the whole localization process or the localization error is larger than 100 m.

As shown in Table VIII, it can be seen that the proposed terrain aided SLAM achieves the highest average localization accuracies in nine scenarios compared to other four methods. The average RMSE of our method is 0.45 m, compared to ORB-SLAM3 0.75 m, ORB-SLAM2 1.32 m, SURF-SLAM 1.74 m, and SuperPoint-SLAM 1.61 m.

ORB-SLAM2 and ORB-SLAM3 have similar localization performance, while ORB-SLAM3 has a better performance, especially in the aspect of maximum error due to the better place recognition algorithm that provides more midterm matches. However, it should be noticed that ORB-SLAM3 fails to complete localization in scenario 9 because of tracking failure.

Both SURF-SLAM and SuperPoint-SLAM failed twice in nine scenarios. These failures due to the reason that UAV images from 01 to 03 scenarios are generated from real Mars satellite images and thus contain many textureless areas, which is difficult for feature extraction and tracking in SURF-SLAM and SuperPoint-SLAM.

Since the total length of flight trajectories of the UAVs in the nine scenes is 1069.9 m, the localization error of our proposed method is 0.042%. This result confirms that the LBA is able to fuse the geo-referencing and visual odometry into an optimized result. Finally, the processing speed of the proposed method is 12 frames/s, which ensures the real-time performance.

VI. EXPERIMENT RESULTS

In this article, a visual SLAM method based on terrain aided geo-referencing is proposed for planetary UAV optical navigation. Based on the mathematical derivation to prove the robustness of PC to multimodal image registration, a PC-based geo-referencing method is proposed to automatically localize the UAV image sequence on the reference terrain model. The proposed optical navigation approach includes three main steps: PC for terrain-based geo-referencing, visual odometry based on frame-to-frame feature tracking, and keyframe pose refinement based on LBA. Image matching experiments between DEM and optical images and optical navigation experiments using simulated Mars UAV sequences from different terrain models have been conducted to rigorously assess the capability of the optical navigation method. The results demonstrated that compared to the state-of-the-art image matching methods, the proposed terrain aided navigation method is able to robustly recover the UAV trajectories, regardless of the large appearance difference between UAV and reference DEM data.

It should be noted that our navigation method currently only relies on monocular camera, and future work can be carried on using multisensor, such as stereo camera and IMU. Another promising direction is to involve multiple reference maps, such as orthoimages, vector maps, and semantic maps, for robust geo-referencing. Matching failures are hardly to avoid when reference maps and onboard UAV images come from different sources, it would be beneficial to broaden the sources of reference maps to achieve more robust navigation results.

The work presented in this article laid the theoretical foundation for terrain aided UAV optical navigation. It can also be applied in vision-based UAV navigation on the Earth

TABLE VIII
POSITIONING ACCURACY OF THE VISUAL ODOMETER WITH TERRAIN ASSISTANCE IN EACH SCENE

Scene	method	max/m	min/m	mean/m	rmse/m	std/m
01	ORB-SLAM2	5.69	0.08	1.42	1.60	0.73
	ORB-SLAM3	2.51	0.05	0.57	0.69	0.39
	SURF-SLAM	-	-	-	-	-
	SuperPoint-SLAM	-	-	-	-	-
	Our Method	4.64	0.02	0.45	0.60	0.39
02	ORB-SLAM2	2.16	0.16	0.87	0.92	0.29
	ORB-SLAM3	2.95	0.04	0.49	0.57	0.28
	SURF-SLAM	3.25	0.18	1.17	1.28	0.52
	SuperPoint-SLAM	-	-	-	-	-
	Our Method	2.03	0.01	0.30	0.37	0.21
03	ORB-SLAM2	1.62	0.04	0.67	0.71	0.25
	ORB-SLAM3	1.35	0.04	0.41	0.45	0.19
	SURF-SLAM	-	-	-	-	-
	SuperPoint-SLAM	-	-	-	-	-
	Our Method	1.48	0.01	0.26	0.31	0.17
04	ORB-SLAM2	3.52	0.19	1.02	1.15	0.54
	ORB-SLAM3	3.02	0.24	1.18	1.29	0.51
	SURF-SLAM	4.78	0.31	1.07	1.15	0.42
	SuperPoint-SLAM	4.63	0.33	2.27	2.34	0.58
	Our Method	2.30	0.03	0.26	0.34	0.21
05	ORB-SLAM2	1.78	0.01	0.45	0.52	0.25
	ORB-SLAM3	1.64	0.03	0.41	0.47	0.23
	SURF-SLAM	1.85	0.12	0.76	0.81	0.27
	SuperPoint-SLAM	3.10	0.10	1.14	1.32	0.66
	Our Method	0.76	0.01	0.22	0.26	0.13
06	ORB-SLAM2	2.02	0.11	0.89	0.95	0.31
	ORB-SLAM3	1.52	0.10	0.68	0.72	0.26
	SURF-SLAM	4.10	0.44	1.32	1.41	0.49
	SuperPoint-SLAM	3.95	0.10	1.02	1.14	0.50
	Our Method	3.56	0.02	0.30	0.42	0.29
07	ORB-SLAM2	3.09	0.21	0.93	1.00	0.35
	ORB-SLAM3	2.43	0.15	0.83	0.88	0.29
	SURF-SLAM	4.93	0.11	1.16	1.29	0.57
	SuperPoint-SLAM	6.57	0.20	1.46	1.67	0.80
	Our Method	3.89	0.02	0.50	0.69	0.48
08	ORB-SLAM2	5.61	0.35	1.92	2.13	0.93
	ORB-SLAM3	2.70	0.08	0.84	0.91	0.35
	SURF-SLAM	8.06	1.49	3.29	3.54	1.31
	SuperPoint-SLAM	3.80	0.34	1.90	2.02	0.68
	Our Method	2.86	0.007	0.29	0.38	0.24
09	ORB-SLAM2	5.83	0.77	2.72	2.89	0.98
	ORB-SLAM3	-	-	-	-	-
	SURF-SLAM	5.53	0.48	2.45	2.69	1.12
	SuperPoint-SLAM	3.40	0.15	1.02	1.19	0.60
	Our Method	5.77	0.05	0.44	0.68	0.53

when the GPS signal is not accurate or unavailable. It can also be applied in the vision-based landing process on an asteroid or

other planet surface when a terrain model is available. Future work will be carried out using the real UAV data captured from

the desert area, which has a similar texture to the planetary surface.

REFERENCES

- [1] L. Pllice, B. Lau, G. Pisanich, and L. A. Young, "Biological inspired behavioral strategies for autonomous aerial explorers on Mars," in *Proc. IEEE Aerosp. Conf.*, Mar. 2003, p. 304.
- [2] J. B. Balaram, "Mars helicopter," Tech. Rep., 2020.
- [3] P. Turtle *et al.*, "Dragonfly: Exploring Titan's prebiotic organic chemistry and habitability," *Icarus*, vol. 243, pp. 191–207, 2007.
- [4] L. Matthies *et al.*, "Computer vision on Mars," *Int. J. Comput. Vis.*, vol. 75, no. 1, pp. 67–92, Jul. 2007.
- [5] K. Janschek, V. Janschek, and M. Beck, "Performance analysis for visual planetary landing navigation using optical flow and DEM matching," in *Proc. AIAA Guid., Navig., Control Conf. Exhib.*, Aug. 2006, p. 6706.
- [6] Y. Cheng, M. Maimone, and L. Matthies, "Visual odometry on the Mars exploration rovers," in *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, vol. 1, Oct. 2005, pp. 903–910.
- [7] S. Chiodini *et al.*, "Mars rovers localization by matching local horizon to surface digital elevation models," in *Proc. IEEE Int. Workshop Metrology Aerosp. (MetroAeroSpace)*, Jun. 2017, pp. 374–379.
- [8] F. Cozman, E. Krotkov, and C. Guestrin, "Outdoor visual position estimation for planetary rovers," *Auto. Robots*, vol. 9, no. 2, pp. 135–150, 2000.
- [9] E. Fang, "Terrain relative navigation for lunar polar roving: Exploiting geometry, shadows, and planning," Carnegie Mellon Univ., Pittsburgh, PA, USA, Tech. Rep., 2020.
- [10] A. V. Nefian *et al.*, "Planetary rover localization within orbital maps," in *Proc. IEEE Int. Conf. Image Process. (ICIP)*, Oct. 2014, pp. 1628–1632.
- [11] T. Miso, T. Hashimoto, and K. Ninomiya, "Optical guidance for autonomous landing of spacecraft," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 35, no. 2, pp. 459–473, Apr. 1999.
- [12] A. E. Johnson, Y. Cheng, and L. H. Matthies, "Machine vision for autonomous small body navigation," in *Proc. IEEE Aerosp. Conf.*, Mar. 2000, pp. 661–671.
- [13] Y. Cheng, A. Johnson, and L. Matthies, "MER-DIMES: A planetary landing application of computer vision," in *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2005, pp. 806–813.
- [14] Z. Liu *et al.*, "Remote sensing mapping and localization techniques for teleoperation of chang'e-3 rover," *J. Remote Sens.*, vol. 18, pp. 971–980, 2014.
- [15] H. Bay, T. Tuytelaars, and L. V. Gool, "Surf: Speeded up robust features," in *Proc. Eur. Conf. Comput. Vis.* Springer, 2006, pp. 404–417.
- [16] R. Mur-Artal and J. D. Tardós, "ORB-SLAM2: An open-source slam system for monocular, stereo, and RGB-D cameras," *IEEE Trans. Robot.*, vol. 33, no. 5, pp. 1255–1262, Oct. 2017.
- [17] C. Campos, R. Elvira, J. J. G. Rodriguez, J. M. M. Montiel, and J. D. Tardos, "ORB-SLAM3: An accurate open-source library for visual, visual-inertial, and multimap SLAM," *IEEE Trans. Robot.*, vol. 37, no. 6, pp. 1874–1890, Dec. 2021.
- [18] D. DeTone, T. Malisiewicz, and A. Rabinovich, "SuperPoint: Self-supervised interest point detection and description," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2018, pp. 224–236.
- [19] A. Papoulis and S. U. Pillai, *Probability, Random Variables, and Stochastic Processes*. New York, NY, USA: McGraw-Hill, 2002.
- [20] E. Shechtman and M. Irani, "Matching local self-similarities across images and videos," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2007, pp. 1–8.
- [21] A. Chilian and H. Hirschmuller, "Stereo camera based navigation of mobile robots on rough terrain," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Oct. 2009, pp. 4571–4576.
- [22] A. Nefian, L. Edwards, D. Lees, L. Keely, T. Parker, and M. Malin, "Automatic rover localization in orbital maps," in *Proc. Lunar Planet. Sci. Conf.*, 2017, p. 2374.
- [23] M. Balcı and H. Foroosh, "Subpixel estimation of shifts directly in the Fourier domain," *IEEE Trans. Image Process.*, vol. 15, no. 7, pp. 1965–1972, Jul. 2006.
- [24] W. S. Hoge, "A subspace identification extension to the phase correlation method [MRI application]," *IEEE Trans. Med. Imag.*, vol. 22, no. 2, pp. 277–280, Feb. 2003.
- [25] J. G. Liu and H. Yan, "Phase correlation pixel-to-pixel image co-registration based on optical flow and median shift propagation," *Int. J. Remote Sens.*, vol. 29, no. 20, pp. 5943–5956, Oct. 2008.
- [26] V. Argyriou and T. Vlachos, "Using gradient correlation for sub-pixel motion estimation of video sequences," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process.*, May 2004, p. 329.
- [27] X. Wan, J. G. Liu, and H. Yan, "The illumination robustness of phase correlation for image alignment," *IEEE Trans. Geosci. Remote Sens.*, vol. 53, no. 10, pp. 5746–5759, Oct. 2015.
- [28] X. Wan, J. Liu, S. Li, J. Dawson, and H. Yan, "An illumination-invariant change detection method based on disparity saliency map for multitemporal optical remotely sensed images," *IEEE Trans. Geosci. Remote Sens.*, vol. 57, no. 3, pp. 1311–1324, Mar. 2018.
- [29] X. Wan, J. Liu, H. Yan, and G. L. Morgan, "Illumination-invariant image matching for autonomous UAV localisation based on optical sensing," *ISPRS J. Photogramm. Remote Sens.*, vol. 119, pp. 198–213, Sep. 2016.
- [30] X. Wan, C. Wang, and S. Li, "The extension of phase correlation to image perspective distortions based on particle swarm optimization," *Sensors*, vol. 19, no. 14, p. 3117, Jul. 2019.
- [31] P. Kube and A. Pentland, "On the imaging of fractal surfaces," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. PAMI-10, no. 5, pp. 704–707, Sep. 1988.
- [32] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM: A versatile and accurate monocular SLAM system," *IEEE Trans. Robot.*, vol. 31, no. 5, pp. 1147–1163, Oct. 2015.
- [33] B. Triggs, P. F. McLauchlan, R. I. Hartley, and A. W. Fitzgibbon, "Bundle adjustment—A modern synthesis," in *Proc. Int. Workshop Vis. Algorithms*. Springer, 1999, pp. 298–372.



Xue Wan received the B.Eng. and M.Eng. degrees from the School of Remote Sensing and Information Engineering, Wuhan University, Wuhan, China, in 2010 and 2012, respectively, and the Ph.D. degree from Imperial College London, London, U.K., in 2015.

She is currently a Professor with the Technology and Engineering Center for Space Utilization, Chinese Academy of Sciences, Beijing, China. Her current research interests are remotely sensed image matching, change detection, vision-based navigation, and 3-D reconstruction.



Yuanbin Shao received bachelor's degree from the School of Astronautics, Nanjing University of Aeronautics and Astronautics, Nanjing, China, in 2020. He is currently pursuing the master's degree with the School of Aeronautics and Astronautics, University of Chinese Academy of Sciences, Beijing, China.

His research interest is the visual localization and navigation of intelligent robots.



Shengyang Zhang is currently pursuing the bachelor's degree in computer science and technology with the School of Information Science and Technology, Dalian Maritime University, Dalian, China.

His current research interests are image matching and vision-based navigation.



Shengyang Li received the M.S. degree in computer science and technology from the Shandong University of Science and Technology, Qingdao, China, and the Ph.D. degree from the Institute of Remote Sensing Applications, Chinese Academy of Sciences, Beijing, China, in 2006.

He is currently a Professor with the Technology and Engineering Center for Space Utilization, Chinese Academy of Sciences. His current research activities are ground data processing system technology, target detection and tracking in video satellites, and machine learning in remote sensing image classification and recognition.