# A Robust Pavement Mapping System Based on Normal-Constrained Stereo Visual Odometry

Huaiyang Huang[1*], Rui Fan[1*], *Member, IEEE*, Yilong Zhu[1],
Ming Liu[1], *Senior Member, IEEE*, Ioannis Pitas[2], *Fellow, IEEE*
[1]Robotics and Multi-Perception Laboratory, Robotics Institute,
the Hong Kong University of Science and Technology, Hong Kong SAR, China.
[2]Department of Informatics, Aristotle University of Thessaloniki, Thessaloniki, Greece.
Emails: {hhuangat, eeruifan, yzhubr, eelium}@ust.hk, pitas@csd.auth.gr

*Abstract*— Pavement condition is crucial for civil infrastructure maintenance. This task usually requires efficient road damage localization, which can be accomplished by the visual odometry system embedded in unmanned aerial vehicles (UAVs). However, the state-of-the-art visual odometry and mapping methods suffer from large drift under the degeneration of the scene structure. To alleviate this issue, we integrate normal constraints into the visual odometry process, which greatly helps to avoid large drift. By parameterizing the normal vector on the tangential plane, the normal factors are coupled with traditional reprojection factors in the pose optimization procedure. The experimental results demonstrate the effectiveness of the proposed system. The overall absolute trajectory error is improved by approximately 20%, which indicates that the estimated trajectory is much more accurate than that obtained using other state-of-the-art methods.

## I. INTRODUCTION

Inspecting and repairing pavement damage is an essential task for public infrastructure maintenance [1]. Manual visual inspection is still the main form of pavement damage inspection [2], which is, however, very labor-intensive and time-consuming [3]. To overcome these drawbacks, more attention is being paid toward developing an automated pavement inspection system [4]. However, these solutions are still not robust and precise enough [5]. Therefore, how to develop an accurate and efficient pavement damage inspection system is still an open problem.

Recent advances in airborne technology make efficient pavement inspection a more solvable problem [6]. Among these techniques, visual odometry (VO) and mapping are two essential modules for an automated pavement inspection system deployed on an unmanned aerial vehicle (UAV). VO provides UAV systems with a fundamental capability for real-time pose estimation, especially onboard perception sensors [7], while the mapping module allows map establishment and relocalization, and thus global position labeling for pavement damage [8]. Recently, researchers have successfully established autonomous pavement inspection systems for UAVs. For example, Zhang *et al.* [9] designed a robust photogrammetric mapping system for UAVs, which can recognize different types of pavement damages, e.g.,
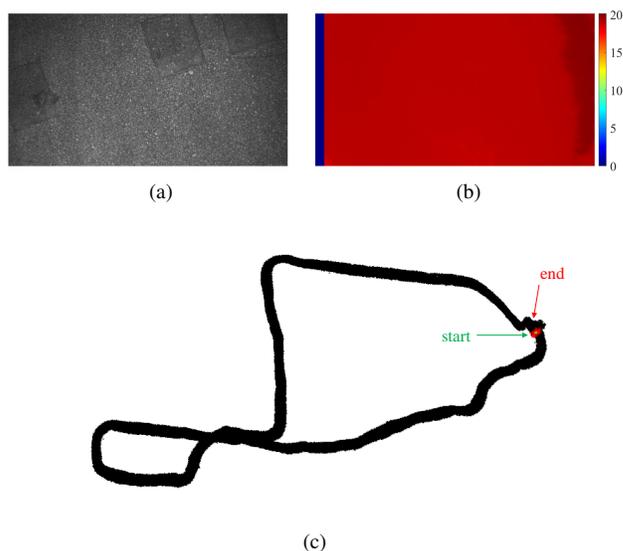
Fig. 1: Pavement mapping result under scene degeneration. An input gray-scale image (a) and the corresponding disparity map (b) for a selected keyframe are shown. (c) Pavement mapping result produced by the proposed visual odometry and mapping system.

cracks and potholes, from RGB images. Furthermore, Fan *et al.* [6] proposed an efficient binocular system that is capable of effectively distinguishing road damage from a transformed disparity map [10].

Current visual odometry and mapping frameworks have demonstrated their accuracy and robustness on various open-source datasets [8], [11], [12]. However, for these state-of-the-art approaches, structure degeneration of visual measurements usually leads to performance degradation in the context of pavement mapping [13]. An example gray-scale image and its corresponding disparity map are shown in Fig. 1, which shows a near-planar structure and well represents the degeneration issue under this scenario. To alleviate this problem, we integrate normal constraints into camera pose estimation. By explicitly minimizing local normal measurements with the global normal prior, we implement a driftless stereo VO for pavement reconstruction. A sample output of the proposed system is shown in Figure 1.(c).
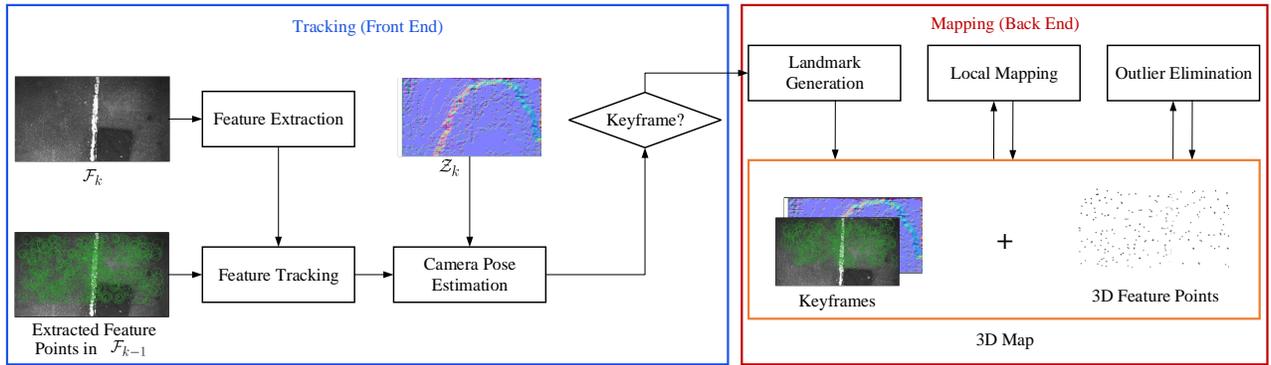
Fig. 2: Framework of the proposed visual odometry system.

The remainder of this paper is structured as follows: Section. I reviews the state-of-the-art visual odometry methods for UAV pose estimation. Section. III details the proposed pavement mapping system. Section. IV presents the experimental results and discusses the performance of our5 system both qualitatively and quantitatively. Finally, Section. V summarizes the paper and provides recommendations for future work.

## II. STATE OF THE ART

The state-of-the-art VO or simultaneous localization and mapping (SLAM) systems can be categorized as indirect [13], [14], direct [11] or hybrid [15] methods. Recent progress in these methods provides fundamental building blocks for onboard UAV pose estimation both theoretically and technically. While monocular SLAM algorithms can not recover reliable scale information [16]. Recent approaches generally resort to different sensor configurations, including stereo vision or visual-inertial sensors.

Some consider introducing stereo constraints for scale recovery. For instance, Cvišic et al. [12] proposed a stereo SLAM framework that is highly efficient in computational demand. To save computational resources, Sun et al. [17] introduce stereo constraints into a filter-based visual-inertial odometry framework, which was named as S-MSCKF. Raul et al. proposed ORB-SLAM2 [8], a versatile visual SLAM framework equipped with sparse mapping, loop-closure and relocalization ability.

Another common strategy is to fuse visual state estimation with an inertial measurement unit (IMU). This research track tackles the problem based on both filtering and optimization methods. Weiss et al. [18] introduced a loosely-coupled filter to recover absolute scale with the aid of an IMU. 6-DoF poses initially estimated by PTAM [14] are fused with the IMU measurements. For tightly-coupled filtering, Bloesch et al. [19] extracted multi-level patch features along with 3D landmarks in the camera tracking procedure. Then camera poses are estimated by a standard extended Kalman Filter. Leutenegger et al. [20] established a sliding window-based optimization framework with keyframe selection. To optimize camera poses, they formulated a cost function combining both visual reprojection error and inertial error terms. Forster et al. [21] proposed a pre-integration strategy

to bootstrap visual-inertial odometry, while Qin et al. [7] introduce a loosely-coupled fusion procedure to initialize parameters including scale and bias. With IMU measurements pre-integrated, a tightly-coupled back-end jointly optimize camera poses along with other parameters.

## III. METHODOLOGY

### A. Notation

Throughout the paper, we denote the image collected at the $k$-th time as $I_k$ and the corresponding frame as $\mathcal{F}_k$. The world coordinate system $\mathcal{F}_w$ is identical to the first camera coordinated system $\mathcal{F}_0$.

For $I_k$, the rigid transform $\mathbf{T}_k \in \mathbf{SE}(3)$ maps a 3D landmark $\mathbf{p} \in \mathbb{R}^3$ to the camera frame using [22]:

$$\mathbf{p}_k^c = \mathbf{R}_k \mathbf{p}_i + \mathbf{t}_k, \tag{1}$$

where $\mathbf{T}_k = [\mathbf{R}_k | \mathbf{t}_k]$, $\mathbf{T}_k \in \mathbf{SE}(3)$. $\mathbf{R}_k$ and $\mathbf{t}_k$ are the rotational and translational components of $\mathbf{T}_k$, respectively. Accordingly, $\mathbf{p}^c$ denotes a 3D point in $\mathcal{F}_k$.

We use $\pi : \mathbb{R}^3 \to \mathbb{R}^2$ to denote the projection function: $\mathbf{u} = \pi(\mathbf{p}^c)$, where $\mathbf{u}$ is a pixel in the image coordinate system (ICS). $\pi$ is defined as [23]:

$$\pi\left(\begin{bmatrix} X \\ Y \\ Z \end{bmatrix}\right) = \begin{bmatrix} f_x \frac{X}{Z} + c_x \\ f_y \frac{Y}{Z} + c_y \\ f_x \frac{X-b}{Z} + c_x \end{bmatrix}, \tag{2}$$

where $(f_x, f_y)$ represents the focal lengths (in pixels) in the horizontal and vertical directions, $(c_x, c_y)$ denotes the principal (in pixels), and $b$ is the baseline of the stereo rig .

The update of a camera pose is parameterized as an incremental twist $\boldsymbol{\xi} \in \mathfrak{se}(3)$. We use a left-multiplicative formulation $\oplus : \mathfrak{se}(3) \times \mathbf{SE}(3) \to \mathbf{SE}(3)$ to denote the update of $\mathbf{T}_k$, which is denoted as

$$\boldsymbol{\xi} \oplus \mathbf{T}_k := \exp(\boldsymbol{\xi}^\wedge) \cdot \mathbf{T}_k. \tag{3}$$

### B. Overview

An overview of the proposed system is shown in Fig. 2. In the preprocessing stage, for each input frame $I_k$, we compute a dense disparity map and normal map. Then the relevant information of $I_k$ is delivered to the VO module. Firstly, the front-end tracker extracts features from accelerated segment test (FAST) [24] feature points and computes their corresponding Oriented FAST and Rotated BRIEF (ORB)
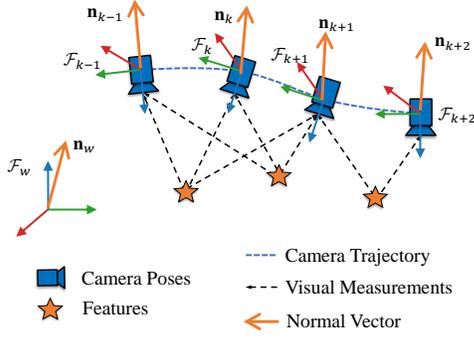
Fig. 3: An illustration for camera pose estimation.

descriptors [25]. With feature association across consequential frames, the current camera pose is estimated through a Perspective-n-Point (PnP) scheme [26]. After initial tracking, the current frame with corresponding matches is delivered to the mapping module, if it is selected as a keyframe. The back-end optimizer jointly estimates camera poses and 3D positions of landmarks.

### C. Pose Estimation with Normal Constraints

This section explains how we introduce the normal constrains into camera pose optimization. Fig. 3 illustrates the Bundle Adjustment (BA) procedure. We combine two types of factors for camera pose estimation, visual residuals $\mathbf{e}_{i,k}^{\text{repro}}$ and normal residuals $\mathbf{e}_k^{\text{normal}}$.

We use the reprojection error as the visual constraint. The residual term defined on the $i$-th landmark and the $k$-th keyframe is defined as:

$$\mathbf{e}_{i,k}^{\text{repro}} = \pi\left(\mathbf{R}_k \mathbf{p}_i + \mathbf{t}_k\right) - \mathbf{u}_{i,k}, \qquad (4)$$

where $\mathbf{u}_{i,k}$ is the pixel coordinates of the feature associated with $\mathbf{p}_i$ in the $k$-th keyframe. Generally, the visual constraint measures the distance between the projected position and the observed position.

Although the traditional BA pipeline is able to estimate camera poses, large drift (especially rotation) is observed in the experiments in the context of pavement mapping. We attribute the failure to strong scene degeneration, as shown in Fig. 1. To tackle this issue, we integrate the normal constraints into the pose optimization procedure. We estimate the normal of the local structure based on the assumption that the 3D geometry is composed of a set of planar surfaces. Then, by minimizing the residual between global and local surface normals, these factors contribute to the estimation of the rotational component of the camera poses.

The normal constraints are introduced to provide additional observations under the degenerated scenario. We leverage the normal-based regulation based on two observations: 1). the pavements in the man-made world typically have a near-planar characteristic; 2). the structure observed in a single frame is a planar surface, which can also be calculated with a closed-form formulation. We estimate the normal $\mathbf{n}_k$ in each frame. Inspired by [7], we define this residual term on the tangential plane orthogonal to $\mathbf{n}_k$ under $\mathcal{F}_k$ as:

$$\mathbf{e}_k^{\text{normal}} = \mathbf{B}_k\left(\mathbf{R}_k \frac{\mathbf{n}_w}{\|\mathbf{n}_w\|} - \mathbf{n}_k\right), \qquad (5)$$
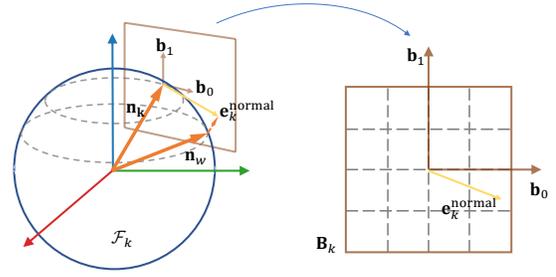


Fig. 4: The parameterization of normal vector and corresponding residual.

where $\mathbf{B}_k = [\mathbf{b}_{k0}, \mathbf{b}_{k1}]^T, \mathbf{B}_k \in \mathbb{R}^{2 \times 3}$ consists of two base vectors $\mathbf{b}_{k0}, \mathbf{b}_{k1}$ of the tangential plane. Fig. 4 illustrates the definition of the residual terms related to the normal constraints. From the geometrical perspective, it projects the difference of the global normal and local normal in $\mathcal{F}_k$ onto the tangential plane, yielding a residual term in $\mathbb{R}^2$. In practice, this avoids overparameterization of normal vectors that can mislead the optimization process. Additionally, unlike in [7], the optimal normal $\mathbf{n}_k$ observed in $\mathcal{F}_k$ is constant in this formulation. Therefore, there is no need to recompute $\mathbf{B}_k$ in each iteration, which accelerates the process of optimization.

To generate the two base vectors $\mathbf{b}_{k0}, \mathbf{b}_{k1}$, we arbitrarily assign a unit vector $\mathbf{v}$ that is not parallel to $\mathbf{n}_k$. Then $\mathbf{b}_{k0}, \mathbf{b}_{k1}$ can be solved by:

$$\mathbf{b}_{k0} = \frac{\mathbf{n}_k \times \mathbf{v}}{\|\mathbf{n}_k \times \mathbf{v}\|}, \qquad \mathbf{b}_{k1} = \frac{\mathbf{n}_k \times \mathbf{b}_{k0}}{\|\mathbf{n}_k \times \mathbf{b}_{k0}\|}. \qquad (6)$$

The proposed residual definition is more efficient, as there is no need to recompute the bases of the tangential plane when the camera pose is updated.

To estimate the camera pose $\mathbf{T}_k = [\mathbf{R}_k, \mathbf{t}_k]$ of the $k$-th frame, we sum over all the valid factors. Then the camera pose is given by

$$\mathbf{T}_k = \arg\min_{\mathbf{T}_k} \sum_{i \in \mathcal{P}_k} w_{i,k} \|\mathbf{e}_{i,k}^{\text{repro}}\|_\gamma + \lambda \cdot \|\mathbf{e}_k^{\text{normal}}\|_\gamma, \qquad (7)$$

where $\mathcal{P}_k$ is the set of landmarks matched successfully in the current frame. $\|\cdot\|_\gamma$ stands for the robust Huber norm and $w_{i,k}$ represents the optimization weight associated with pixel variance. We use a constant factor $\lambda$ to balance the contribution of different factors to the pose optimization.

### D. Back-end Optimization

The optimization back-end uses a similar objective function as Eq. 7. The parameters to be optimized are denoted as $\mathcal{X} = \{\mathbf{p}_i, \mathbf{T}_k | i \in \mathcal{P}, k \in \mathcal{T}\}$, where $\mathcal{P}$ and $\mathcal{T}$ stores the indices of the keyframes and the landmarks in the optimization window, respectively. We have $\mathcal{P} = \bigcup_{k \in \mathcal{T}} \mathcal{P}_k$ Therefore, a full BA is formulated as:

$$\mathcal{X} = \arg\min_{\mathcal{X}} \sum_{k \in \mathcal{T}} \left[ \lambda \cdot \|\mathbf{e}_k^{\text{normal}}\|_\gamma + \sum_{i \in \mathcal{P}_k} w_{i,k} \|\mathbf{e}_{i,k}^{\text{repro}}\|_\gamma \right]. \qquad (8)$$

Note that, as in [8], the poses of keyframes that do not directly connect to the current keyframe in the covisibility

Fig. 5: Experimental set-up for acquiring stereo road images.



(a)          (b)

Fig. 6: Comparison between the trajectory estimated by ORB-SLAM2 and ours in the z-axis.

graph are fixed during the optimization. Like in [8], we detect and reject outliers by $\mathcal{X}^2$-test. Assuming one-pixel variance for every feature, we have $th_{\text{repro}} = 7.815$. For the factor $e_{i,j} > th_{\text{repro}}$, the corresponding feature $\mathbf{p}_i$ and related observations will be rejected by the mapping module.

To initialize the global normal estimation, we add the global normal into the optimizable parameter set $\mathcal{X}$, yielding an augmented set $\mathcal{X}_{\text{init}} = \mathcal{X} \bigcup \{\mathbf{n}_w\}$. We set a window size $\Delta$ and keep $\mathbf{n}_w$ in the optimizable parameter set until $\Delta$ keyframes have passed to the back-end. Then $\mathbf{n}_w$ is fixed in every subsequent optimization.

## IV. EXPERIMENTAL RESULTS

In this section, we present both qualitative and quantitative experimental results of the proposed normal-constrained stereo VO. We first describe the experimental set-up and then compare our system with ORB-SLAM2 [8] to demonstrate the effectiveness.

### A. Experimental Set-Up

In the experiments, an Intel RealSense stereo camera D435i[1] was mounted on a DJI Matrice 100 drone[2] to acquire stereo road images. The maximum take-off weight of the drone is 3.6 kg. The stereo camera can capture stereo images with resolution of $1696 \times 480$ at a speed of 30 fps. An NVIDIA Jetson TX2 embedded system[3] was utilized to save the captured stereo images. An illustration of the experimental set-up is shown in Fig. 5. Using the above experimental set-up, four datasets including 26372 stereo image pairs were created. The stereo camera was calibrated manually using the stereo calibration toolbox from MATLAB R2019a. The resolution of the rectified reference and target images is $824 \times 449$. Also, we mounted a DJI N3 GPS[4] sensor on the UAV to acquire the flight trajectory ground truth. The GPS precision is $\pm 1$ m. The datasets will be publicly available at `https://www.ram-lab.com/`.

### B. Qualitative Evaluations

Fig. 6 compares the proposed system against ORB-SLAM2 [8]. As shown in Fig. 6, the odoemtry drift is reduced significantly in the proposed system. This can be attributed to introducing the normal constraints, so the rotation
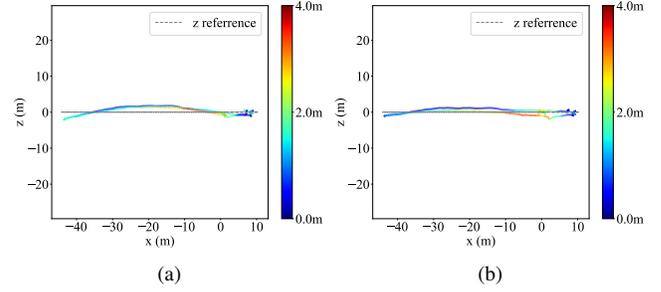
estimations converge to a better minimum, which benefits the translation estimations in the z-axis.

### C. Quantitative Evaluations

In the experiments, we compare the proposed VO algorithm with ORB-SLAM2 [8], a state-of-the-art stereo visual odometry system. We compute the absolute trajectory error (ATE) $e^{\text{ATE}}$ and relative distance error (RDE) $e^{\text{RDE}}$ between the estimated and ground truth trajectories as the metrics for the evaluations. $e_i^{\text{ATE}}$ and $e_i^{\text{RDE}}$, the ATE and RDE of the $i$-th input frame, are similar to the metrics in [27]:

$$e_i^{\text{ATE}} = \Pi(\mathbf{T}_i^{-1}\mathbf{S}\mathbf{T}_i^{\text{gt}}), \tag{9}$$

$$e_i^{\text{RDE}} = |\Pi\left(\mathbf{T}_i^{-1}\mathbf{T}_{i+\Delta}\right)\|_2 - \|\Pi\left((\mathbf{T}_i^{\text{gt}})^{-1}\mathbf{T}_{i+\Delta}^{\text{gt}}\right)\|_2|, \tag{10}$$

where $\mathbf{T}_i$ represents the $i$-th estimated trajectory, $\mathbf{T}_i^{\text{gt}}$ represents the $i$-th ground truth trajectory, $\mathbf{S}$ denotes the rigid-body transformation from $\mathbf{T}_i^{\text{gt}}$ to $\mathbf{T}_i$, $\Pi$ is a function to extract the translation components in the $x$- and $y$-axes, and $\Delta = 20$ is set to measure the RDE. To quantify the accuracy of the estimated trajectory, we compute the mean, median, root mean square error (RMSE) and standard deviation (SD) of both the ATE and RDE.

The quantitative results are shown in Tab. I and Tab. II, respectively. Although ORB-SLAM2 achieves an accurate and consistent camera pose estimation result throughout the experiments, we further improve the accuracy significantly. According to the comparison, the proposed system generally outperforms ORB SLAM2 and the total ATE is improved . Note that in Dataset 2, large drift in the z-axis leads to a large ATE and variance of ORB-SLAM2. In contrast, the ATE of our method on the same sequence is consistent with the others, which indicates that our method successfully alleviates the VO drift in the z-axis. Furthermore, the mean and median of the RDE, emphasizing the drift of estimations, of the proposed system are lower than those of ORB-SLAM2. This prove the effectiveness of the normal constraints in bounding the drift.

To qualify the performance of the proposed VO system, we align the estimated trajectory with its ground truth, as shown in Fig. 7. Our method is better aligned with the ground truth trajectories than ORB-SLAM2.

---

[1] https://click.intel.com/intel-realsense-depth-camera-d435i-imu.html
[2] https://www.dji.com/uk/matrice100
[3] https://developer.nvidia.com/embedded/buy/jetson-tx2
[4] https://www.dji.com/hk-en

TABLE I: ATE of the estimated UAV flight trajectory.

| | ORB-SLAM2 | | | | Ours | | | |
|---|---|---|---|---|---|---|---|---|
| Dataset | Mean (m) | Median (m) | RMSE (m) | SD (m) | Mean (m) | Median (m) | RMSE (m) | SD (m) |
| Dataset 1 | 1.657 | 1.612 | 1.790 | 0.677 | 1.462 | 1.268 | 1.657 | 0.779 |
| Dataset 2 | 4.693 | 4.575 | 5.289 | 2.436 | 2.396 | 2.311 | 2.632 | 1.090 |
| Total | 3.175 | 3.094 | 3.540 | 1.557 | 1.929 | 1.790 | 2.145 | 0.935 |

TABLE II: RDE of the estimated UAV flight trajectory.

| | ORB-SLAM2 | | | | Ours | | | |
|---|---|---|---|---|---|---|---|---|
| Dataset | Mean (m) | Median (m) | RMSE (m) | SD (m) | Mean (m) | Median (m) | RMSE (m) | SD (m) |
| Dataset 1 | 0.160 | 0.109 | 0.230 | 0.166 | 0.116 | 0.072 | 0.186 | 0.146 |
| Dataset 2 | 0.159 | 0.103 | 0.244 | 0.185 | 0.119 | 0.071 | 0.205 | 0.167 |
| Total | 0.149 | 0.106 | 0.238 | 0.177 | 0.118 | 0.072 | 0.196 | 0.157 |



(a) ORB-SLAM2 on Dataset1    (b) ORB-SLAM2 on Dataset2    (c) Ours on Dataset1    (d) Ours on Dataset2
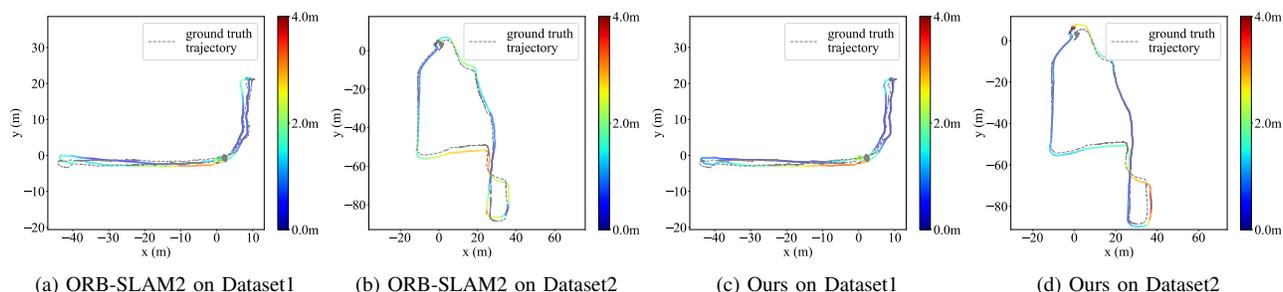
Fig. 7: Comparison between the estimated and ground truth UAV flight trajectory.

TABLE III: Runtime of the tracking and mapping modules.

| Module | Mean (ms) | Median (ms) |
|---|---|---|
| Tracking | 28.7 | 26.5 |
| Mapping | 46.8 | 44.5 |

### D. Timing Results

We measure the runtime of both the tracking and mapping modules, as shown in Tab. III. The proposed system achieves a camera tracking speed of approximately 30 fps when running on an Intel Core i7-8700k CPU. The timing results demonstrate the real-time performance of the proposed normal-constrained stereo SLAM system.

## V. CONCLUSIONS AND FUTURE WORK

In this paper, we presented a pavement mapping system that explicitly introduces ground normal estimations as constraints in the pose optimization and bundle adjustment. We discussed an effective parameterization of normal vectors and corresponding residual term in the optimization. The experimental results showed that our method is drift-less and more accurate than the state-of-the-art ones, which demonstrated the effectiveness of coupling normal constraints with traditional bundle adjustment pipeline.

In the future, we will render the structure assumption become more applicable for pavement mapping. Through minimizing the error between normal observations of single landmarks, the odometry drift might be reduced without explicitly assuming a global constraints. Additionally, with binocular vision only, current implementation is not robust enough under rapid motion (as this may cause motion blur). Therefore, we are planning to leverage inertial measurements and visual-inertial coupling for a more robust visual odometry system. Furthermore, we plan to use our recently published work [28] to provide the UAV with the roll angle information. We believe this can further improve the odometry accuracy.

## REFERENCES

[1] R. Fan, X. Ai, and N. Dahnoun, "Road surface 3D reconstruction based on dense subpixel disparity map estimation," *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 3025–3035, Jun. 2018.

[2] J. S. Miller, W. Y. Bellinger *et al.*, "Distress identification manual for the long-term pavement performance program," United States. Federal Highway Administration. Office of Infrastructure , Tech. Rep., 2014.

[3] R. Fan, U. Ozgunalp, B. Hosking, M. Liu, and I. Pitas, "Pothole detection based on disparity transformation and road surface modeling," *IEEE Transactions on Image Processing*, vol. 29, pp. 897–908, 2020.

[4] R. Fan and M. Liu, "Road damage detection based on unsupervised disparity map segmentation," *IEEE Transactions on Intelligent Transportation Systems*, 2019.

[5] T. Kim and S.-K. Ryu, "Review and analysis of pothole detection methods," *Journal of Emerging Trends in Computing and Information Sciences*, vol. 5, no. 8, pp. 603–608, 2014.

[6] R. Fan, J. Jiao, J. Pan, H. Huang, S. Shen, and M. Liu, "Real-time dense stereo embedded in a uav for road inspection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, Long Beach, CA, USA., Jun. 2019.

[7] T. Qin, P. Li, and S. Shen, "Vins-mono: A robust and versatile monocular visual-inertial state estimator," *IEEE Transactions on Robotics*, vol. 34, no. 4, pp. 1004–1020, 2018.

[8] R. Mur-Artal and J. D. Tardós, "Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras," *IEEE Transactions on Robotics*, vol. 33, no. 5, pp. 1255–1262, 2017.

[9] C. Zhang, "An uav-based photogrammetric mapping system for road condition assessment," *Int. Arch. Photogramm. Remote Sens. Spatial Inf. Sci*, vol. 37, pp. 627–632, 2008.

[10] R. Fan, M. J. Bocus, and N. Dahnoun, "A novel disparity transformation algorithm for road segmentation," *Information Processing Letters*, vol. 140, pp. 18–24, 2018.

[11] J. Engel, V. Koltun, and D. Cremers, "Direct sparse odometry," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 3, pp. 611–625, 2018.

[12] I. Cvišic, J. Cesic, I. Markovic, and I. Petrovic, "Soft-slam: Computationally efficient stereo visual slam for autonomous uavs," *Journal of field robotics*, 2017.

[13] R. Mur-Artal, J. M. M. Montiel, and J. D. Tardos, "Orb-slam: a versatile and accurate monocular slam system," *IEEE Transactions on Robotics*, vol. 31, no. 5, pp. 1147–1163, 2015.

[14] G. Klein and D. Murray, "Parallel tracking and mapping for small ar workspaces," in *Proceedings of the 2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality*. IEEE Computer Society, 2007, pp. 1–10.

[15] C. Forster, M. Pizzoli, and D. Scaramuzza, "Svo: Fast semi-direct monocular visual odometry," in *2014 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2014, pp. 15–22.

[16] R. Fan, J. Jiao, H. Ye, Y. Yu, I. Pitas, and M. Liu, "Key ingredients of self-driving cars," *arXiv:1906.02939*, 2019.

[17] K. Sun, K. Mohta, B. Pfrommer, M. Watterson, S. Liu, Y. Mulgaonkar, C. J. Taylor, and V. Kumar, "Robust stereo visual inertial odometry for fast autonomous flight," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 965–972, 2018.

[18] S. Weiss, M. W. Achtelik, S. Lynen, M. C. Achtelik, L. Kneip, M. Chli, and R. Siegwart, "Monocular vision for long-term micro aerial vehicle state estimation: A compendium," *Journal of Field Robotics*, vol. 30, no. 5, pp. 803–831, 2013.

[19] M. Bloesch, S. Omari, M. Hutter, and R. Siegwart, "Robust visual inertial odometry using a direct ekf-based approach," in *2015 IEEE/RSJ international conference on intelligent robots and systems (IROS)*. IEEE, 2015, pp. 298–304.

[20] S. Leutenegger, P. Furgale, V. Rabaud, M. Chli, K. Konolige, and R. Siegwart, "Keyframe-based visual-inertial slam using nonlinear optimization," *Proceedings of Robotis Science and Systems (RSS) 2013*, 2013.

[21] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-manifold preintegration for real-time visual–inertial odometry," *IEEE Transactions on Robotics*, vol. 33, no. 1, pp. 1–21, 2017.

[22] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.

[23] R. Fan, "Real-time computer stereo vision for automotive applications," Ph.D. dissertation, University of Bristol, 2018.

[24] E. Mair, G. D. Hager, D. Burschka, M. Suppa, and G. Hirzinger, "Adaptive and generic corner detection based on the accelerated segment test," in *European conference on Computer vision*. Springer, 2010, pp. 183–196.

[25] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *Proc. Int. Conf. Computer Vision*, Nov. 2011, pp. 2564–2571.

[26] V. Lepetit, F. Moreno-Noguer, and P. Fua, "Epnp: An accurate o (n) solution to the pnp problem," *International journal of computer vision*, vol. 81, no. 2, p. 155, 2009.

[27] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A benchmark for the evaluation of rgb-d slam systems," in *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 573–580.

[28] R. Fan, L. Wang, M. Liu, and I. Pitas, "A robust roll angle estimation algorithm based on gradient descent," *arXiv preprint arXiv:1906.01894*, 2019.