# ROVO: Robust Omnidirectional Visual Odometry for Wide-baseline Wide-FOV Camera Systems

Hochang Seok and Jongwoo Lim*

{`hochangseok, jlim`}`@hanyang.ac.kr`

Department of Computer Science, Hanyang University, Seoul, Korea.

*Abstract*— In this paper we propose a robust visual odometry system for a wide-baseline camera rig with wide field-of-view (FOV) fisheye lenses, which provides full omnidirectional stereo observations of the environment. For more robust and accurate ego-motion estimation we adds three components to the standard VO pipeline, 1) the hybrid projection model for improved feature matching, 2) multi-view P3P RANSAC algorithm for pose estimation, and 3) online update of rig extrinsic parameters. The hybrid projection model combines the perspective and cylindrical projection to maximize the overlap between views and minimize the image distortion that degrades feature matching performance. The multi-view P3P RANSAC algorithm extends the conventional P3P RANSAC to multi-view images so that all feature matches in all views are considered in the inlier counting for robust pose estimation. Finally the online extrinsic calibration is seamlessly integrated in the backend optimization framework so that the changes in camera poses due to shocks or vibrations can be corrected automatically. The proposed system is extensively evaluated with synthetic datasets with ground-truth and real sequences of highly dynamic environment, and its superior performance is demonstrated.

## I. INTRODUCTION

Ego-motion estimation is an essential functionality in achieving autonomous navigation and maneuver of robots. To estimate the robot's pose in the environment the surrounding structures and objects needs to be modeled accurately. Recently many approaches using various sensors including LIDARs, radars, and/or cameras have been developed. Among them the cameras have advantages for their low cost, passive sensing, abundant information on the environment, mechanical robustness, and many more.

Visual odometry (VO), ego-motion estimation using one or more cameras, has been researched for a few decades to accomplish real-time processing, accurate pose estimation, and robustness to the external disturbances. It has been widely applied to many applications including augmented/virtual reality (AR/VR), advanced driver assistance system (ADAS) and autonomous driving. There exist variety of VO systems, for example ORB-SLAM2 [1] is monocular or stereo and using feature points with binary descriptors, or Stereo-DSO [2] uses a stereo camera and edges for the feature. Monocular systems are attractive for their simple hardware configuration, but it has several limitations that the true scale of motion cannot be estimated, and it is mostly use for AR/VR where metric poses are not needed. Multi-view (stereo) VO systems produce more robust and true
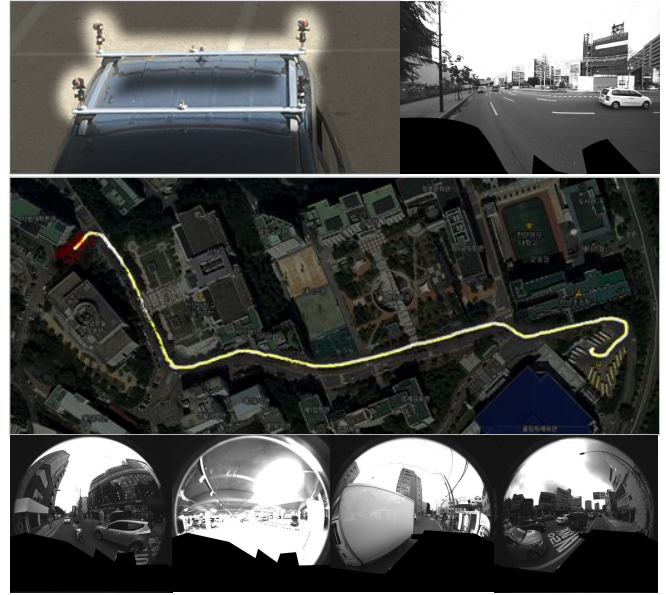


Fig. 1: Top: Our system setup. Four cameras are mounted to the rigid rig on the vehicle rooftop (left). The warped image with our hybrid projection model (right). Middle: the estimated metric trajectory overlaid on the satellite image. Bottom: four wide-FOV input images in challenging real datasets.

metric ego-motion estimates suitable for robot applications, but requires more computational cost. All VO algorithms are vulnerable to excessive dynamic objects, i.e., when the significant portion of the field-of-view (FOV) is covered by moving objects, the pose estimation becomes unstable and incorrect.

In this paper we propose a novel omnidirectional visual odometry system using a multi-camera rig with wide FOV fisheye lenses and wide baseline, to maximize the stability and accuracy of the pose estimation. From four cameras with 220° FOV lenses, it is possible to observe full 360° angle around the robot, and most regions in the environment are visible from more than two views which makes the stereo triangulation possible.

Our wide-baseline wide-FOV setup maximizes the pose accuracy as the stereo resolution is proportional to the baseline length, and it also minimizes the image area occluded by the robot body. However it also poses new challenges. First tracking and matching features become harder because the

viewing directions to the points can be substantially different, and there is large distortion in the periphery of fisheye images where most overlaps between views occur. Second, as the cameras are mounted far apart, the rigid rig assumption is no longer valid. Due to shocks, vibrations, and heat the rig can deform or the cameras can move unexpectedly.

We add three novel components to the VO pipeline to resolve these issues. To improve the feature matching the lens distortion needs to be removed, but as the fisheye lens covers more than $180°$, the image cannot be warped to a single plane. We propose a hybrid projection model that uses two planes in the overlapping regions and a cylinder smoothly connecting them. This projection model enables continuous tracking of feature points in each view and consistent feature descriptors across views.

To estimate the current pose from noisy feature matches, we need to use a RANSAC algorithm with a minimal pose estimator. For the omnidirectional case, one can use minimal solvers with a generalized camera model [3], [4], but we use a simpler approach of computing the pose using P3P [5] from one view and checking the inliers for all feature matches. This multi-view P3P RANSAC effectively and robustly estimates the rig poses in highly dynamic scenes.

Lastly we implement the online extrinsic calibration to deal with unexpected changes of rig-to-camera poses throughout the system execution. Besides the rig deformation by external causes, the initial calibration may not be very accurate due to the size and position of the rig. Online extrinsic calibration built in the local bundle adjustment constantly updates the extrinsic parameters from the tracked features in the current map, and it greatly improves the robustness and accuracy of the system.

For experimental evaluation we render synthetic datasets with ground-truth poses as well as collect challenging real datasets using the omnidirectional rig in Figure 1. Extensive experiments show that the proposed VO system accomplish good performance in all synthetic and real datasets.

Our main contributions can be summarized as follows:

- An effective novel image projection model which allows to find and track feature correspondences between the fisheye cameras in a wide-baseline setup.
- A proposed visual odometry system that uses multiple fisheye cameras with overlapping views operates robustly in highly-dynamic environment using the multi-view P3P RANSAC algorithm and the online extrinsic calibration integrated with the backend local bundle adjustment.
- Extensive evaluation of the proposed algorithms to the synthetic and real datasets verifies the superior performance. All datasets as well as the system implementation will be made public with the paper publication.

## II. RELATED WORKS

In the VO and visual SLAM literature, many different camera configurations have been researched. There are various monocular systems [6]–[8] that are point feature-based, directly optimizing poses with image contents, or
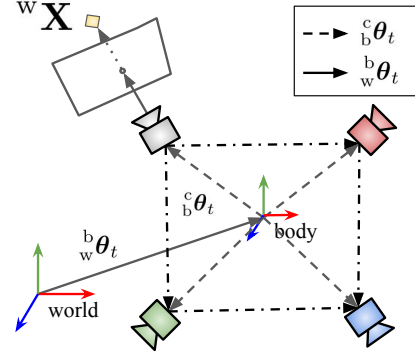


Fig. 2: The world, rig(body), and camera coordinate systems and transformations between them.

hybrid. They show outstanding performance, but due to fundamental limitation of monocular setup, metric poses cannot be estimated. For robotic application, stereo-based systems [1], [2] have been proposed. Another limitation of the conventional systems is small FOV which can make a VO system unstable due to lack of features or existence of dynamic objects. For this practical reason, fisheye camera based methods have been researched recently. Caruso et al. [9] propose a fisheye visual SLAM with direct methods. Liu et al. [10] use a fisheye stereo camera and recover metric scale trajectory. Most recently, Matsuki et al. [11] proposed an omnidirectional visual odometry with the direct sparse method.

For improved environmental awareness and perception capabilities, multi-camera methods also have been studied. [12] present a visual odometry algorithm for a multi-camera system which can observe full surrounding view. They successfully estimate ego-motion with the 2-point algorithm showing the importance of the inter-camera correspondences to recover metric scale. Heng et al. [13] implement a visual SLAM and self-calibration system with at least one calibrated-stereo camera and an arbitrary number of monocular cameras where they have overlapping views with the stereo camera. Recently, a robust multi-camera system using direct methods with plane sweeping stereo is proposed by Liu et al. [14]. Finding correspondences between fisheye images is a challenging and important problem and many researchers devoted efforts in it. Special descriptors [15], [16] are designed to consider the distortion, and Hane et al. [17] and Gao et al. [18] proposed dense matching algorithms for fisheye images.

## III. NOTATION

In this section, we introduce the notations used in this paper. A rigid transformation $\boldsymbol{\theta}$ is parameterized as an axis-angle rotation vector $\mathbf{r}$ and a translation vector $\mathbf{t}$ in $\mathbb{R}^3$. It transforms a 3D point $\mathbf{X}$ to $\boldsymbol{\theta} \star \mathbf{X} = R(\mathbf{r})\mathbf{X} + \mathbf{t}$, where $R(\mathbf{r})$ is the $3 \times 3$ rotation matrix for $\mathbf{r}$. $\star$ also denotes the composition of transformations, and $^{-1}$ is the inverse transformation. As in Figure 2 we use three coordinate systems, world (w), body (b), and camera (c), and when needed the coordinate system is marked on the left of symbols, like $^{\mathrm{w}}_{\mathrm{b}}\boldsymbol{\theta}$ meaning the rigid
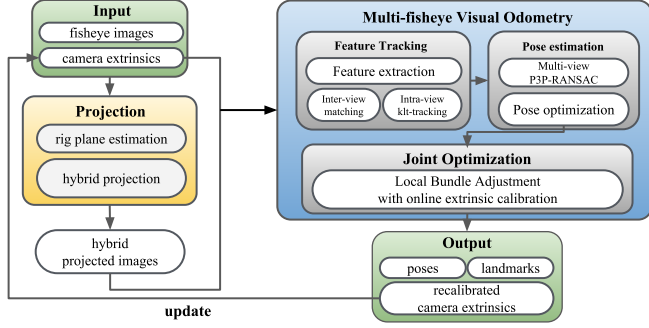
Fig. 3: Overview of the proposed system. The input fisheye images are warped using the hybrid projection model, then the multi-camera VO module tracks and matches the features, computes the current pose, and optimizes the map and the extrinsics. See Section IV for detailed description.

transform from the body to the world coordinate system or $^\mathrm{w}\mathbf{X}$ a point in the world coordinate system. When the time is involved it is denoted as a subscript. For example the camera coordinate of a world point $\mathbf{X}$ at time $t$ can be written as

$$^\mathrm{c}\mathbf{X}_t = {}^\mathrm{c}_\mathrm{b}\boldsymbol{\theta}_t \star {}^\mathrm{b}_\mathrm{w}\boldsymbol{\theta}_t \star {}^\mathrm{w}\mathbf{X}.$$
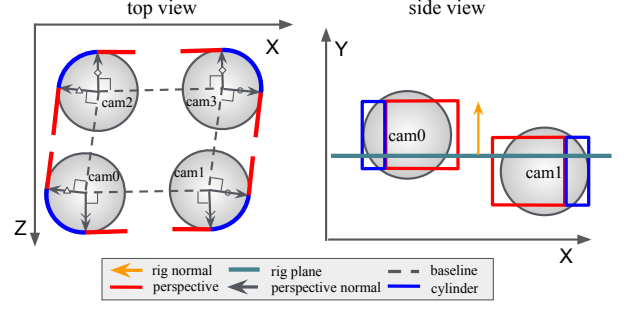
The camera intrinsic calibration determine the mapping between a 3D point $^\mathrm{c}\mathbf{X}$ and a pixel coordinate $\mathbf{x}$ in the image. We denote the projection function $\mathbf{x} = \pi(^\mathrm{c}\mathbf{X}; \phi)$ for a camera intrinsic parameter $\phi$. We use $\pi_0(\cdot)$ for projection onto the unit sphere, $\bar{\mathbf{x}} = \pi_0(\mathbf{X})$, where $\bar{\mathbf{x}}$ is a unit-length ray pointing $\mathbf{X}$, which can also be a feature point in the image.
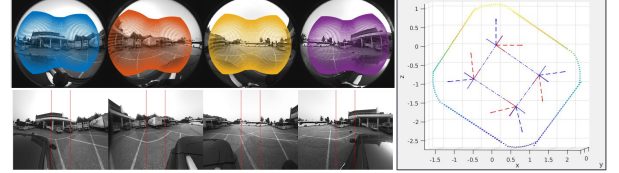
## IV. ALGORITHM

We propose the robust omnidirectional visual odometry (ROVO) system with hybrid projection model, multi-view P3P RANSAC, and online-extrinsic calibration. An overview of the system architecture is show in Figure 3. We assume that the camera intrinsic parameters and the initial extrinsic parameters are known. When input fisheye images arrive, the hybrid projection model warps the input images into perspective-cylindrical images, and the feature detection and tracking modules run on the warped images per view. After the intra-view feature tracking, we perform inter-view feature matching by comparing the feature descriptors in the overlapping regions. Then the multi-view P3P RANSAC algorithm computes the real-scale camera pose, followed by pose optimization which updates the estimated rigid pose with the 3D points fixed. Finally the back-end local bundle adjustment module updates the recent rig trajectory and 3D point locations, as well as the rig extrinsic parameters in a unified framework. For optimization tasks we use Ceres solver [19].

### A. Hybrid Projection Model

Due to wide FOV the original fisheye images contain large amount of information about the environment, but at the same time the periphery of the images is distorted excessively. In our wide-baseline wide-FOV setup there exist



(a) Illustration of the hybrid projection model. The rig plane is best fitting plane given relative position of cameras.



(b) Left: The pixels mapped in the original fisheye images (top) and the warped result image (bottom). Right: the top view of the projection planes and cylinders.

Fig. 4: Hybrid projection model. See Section IV-A

large overlapping area between images, but direct feature matching with the descriptors from such distorted areas does not yield good correspondences in quality and quantity. The local feature descriptors work best when the images are purely perspective, thus we develop a projection model which ensures the overlapping areas between views are as perspective as possible. At the same time, for feature tracking to be successful, the warped image must be continuous and smooth.

To satisfy these conditions, our hybrid projection model has two projection planes at the left and right sides and the cylinder at the center connects the two planes smoothly. Figure 4a shows the hybrid projection model for our rig; for each camera the planes parallel to the baselines of the cameras are connected with cylinders perpendicular to the projection planes. Note that when the camera centers are located close to the rig plane, the proposed method is similar to the stereo rectification and the projected $y$-coordinates of a scene point in the other images should be roughly same.

To build the warped image, we need to find the pixel position $\hat{\mathbf{x}}$ in the original fisheye image corresponding to each pixel $\mathbf{x}$ in the warped image. For the point $\mathbf{x}$, the 3D point $^\mathrm{c}\boldsymbol{\Pi}(\mathbf{x})$ on the projection plane/cylinder can be computed using the plane/cylinder equation, and then its fisheye image coordinate is given as $\hat{\mathbf{x}} = \pi(^\mathrm{c}\boldsymbol{\Pi}(\mathbf{x}); \phi)$. Figure 4b shows the fisheye coordinates for the warped pixels (top) and the warped images (bottom).

### B. Intra- and Inter-view Feature Processing

For each warped image, we perform the intra-view feature processing which is the standard feature detection and tracking. We use the ORB feature detector [20] with minimum distance constraints to neighboring features to ensure that features are extracted uniformly. The existing feature points

are tracked by the KLT tracking algorithm [21].

To improve the pose estimation and mapping quality we need to track the features across views. The inter-view feature processing finds the matching features in the overlapping regions between views and transfer the information when a feature goes out of the FOV in one image. The ORB descriptors are attached to the tracked points, and we use the $k$-nearest neighbor feature matching algorithm to find the feature correspondences. Incorrect matches are filtered according to the policies similar to the stereo matching:

- $y$-distance between the two matched points is small.
- the feature discriptor and orientation difference is small.
- it fulfills epipolar consistency, left-right consistency and positive disparity.
- the zero normalized cross correlation (ZNCC) cost is small.

Finally, we triangulate the matched points to compute their 3D coordinates, which are used in pose estimation. After the feature processing, we obtain the feature-landmark correspondences $\{\{(\bar{\mathbf{x}}_{i_j}, {}^{\mathrm{w}}\mathbf{X}_{i_j})\}_{i_j}\}_j$ for all camera $j$'s. Note the feature locations are convertd to unit-length rays.

### C. Multi-view P3P RANSAC

After the feature processing and triangulation, the current pose of the rig is estimated from the established 2D-3D feature correspondences. Our multi-view P3P RANSAC algorithm extends the monocular P3P RANSAC algorithm [22]. In our RANSAC iterations, one view is selected randomly with Probability Proportional to the Size of match sets (PPS sampling), then the minimum sample set is randomly chosen among the correspondences in the view. From the camera pose candidates estimated by the monocular P3P algorithm, the rig poses are computed and all correspondences in all views are tested for inlier check. The detailed process is shown in Algorithm 1.

PPS sampling choose the cameras with more feature matches more frequently to increase the chance of finding good poses, while all-view inter checking enforces the estimated pose is consistent with all observations. To determine the best pose in the RANSAC loop, we use the reprojection errors of the inlier matches only.

After RANSAC finishes, the estimated rig pose ${}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}$ is optimized by minimizing the reprojection error of all inliers while the 3D points are fixed,

$$\min_{{}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}} \sum_j \sum_{i_j^*} \rho\left(||\bar{\mathbf{x}}_{i_j^*} - \pi_0({}_{\mathrm{b}}^{j}\boldsymbol{\theta} \star {}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta} \star {}^{\mathrm{w}}\mathbf{X}_{i_j^*})||^2\right),$$

where ${}_{\mathrm{b}}^{j}\boldsymbol{\theta}$ is the transformation from the body to the camera $j$'s coordinate system, $i_j^*$'s are the inlier indices, and $\rho$ is the Cauchy loss function which minimizes the influence of outliers.

### D. Online Extrinsic Calibration

To deal with the deformation and motion of the camera during operation, the camera extrinsic parameters are jointly updated in the local bundle adjustment module. For online

---

**Algorithm 1:** multi-view P3P RANSAC algorithm

**Data:** 2D feature locations and 3D landmark correspondences in all views, $\{\{(\bar{\mathbf{x}}_{i_j}, {}^{\mathrm{w}}\mathbf{X}_{i_j})\}_{i_j}\}_j$

**Result:** Rigid body transformation ${}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}^*$

1 **while** $iter < iter_{max}$ **do**
2     Select a camera $j'$ using PPS sampling.
3     Randomly sample 3 pairs in the selected camera.
4     Get camera pose candidates $\{{}_{\mathrm{w}}^{j'}\boldsymbol{\theta}_k\}$ by P3P.
5     Compute the body pose candidates $\{{}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}_k\}$.
6     **for** *each body pose candidate* ${}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}_k$ **do**
7         **for** *each camera $j$* **do**
8             Compute the camera pose ${}_{\mathrm{w}}^{j}\boldsymbol{\theta}_k$.
9             Add reprojection error of all inliers in $j$:
10             $C_k \mathrel{+}= \sum_{i_j}\max(0, \tau_r - ||\bar{\mathbf{x}}_{i_j} - \pi_0({}_{\mathrm{w}}^{j}\boldsymbol{\theta}_k \star {}^{\mathrm{w}}\mathbf{X}_{i_j})||)$
11         **end**
12         **if** $C_k$ *is largest* **then**
13             ${}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}^* \leftarrow {}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}_k$.
14             Update $iter_{max}$ with the new inlier ratio.
15         **end**
16     **end**
17 **end**

---

extrinsic calibration, we add the camera extrinsic parameters $\{{}_{\mathrm{b}}^{\mathrm{c}}\boldsymbol{\theta}\}_{\mathrm{c}}$ into the optimization

$$\min \sum_t \sum_j \sum_{i_j} \omega_{i_j} \rho\left(||\bar{\mathbf{x}}_{i_j,t} - \pi_0({}_{\mathrm{b}}^{j}\boldsymbol{\theta} \star {}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}_t \star {}^{\mathrm{w}}\mathbf{X}_{i_j})||^2\right),$$

where the rig poses $\{{}_{\mathrm{w}}^{\mathrm{b}}\boldsymbol{\theta}_t\}_t$ in the active time window, the landmark positions $\{{}^{\mathrm{w}}\mathbf{X}_i\}_i$, as well as the camera extrinsics $\{{}_{\mathrm{b}}^{j}\boldsymbol{\theta}\}_j$ are optimized to minimize the cost. We give higher weight $\omega_{i_j}$ for the points observed in multiple cameras.

Since the cameras in our system are fixed on a rig, we need to give an extra constraint that the distance between the cameras are constant. Otherwise the metric scale reconstruction is not possible as the rig can grow or shrink over time. The constraints can be written as $||{}_{i}^{j}\mathbf{t}|| = ||{}_{i}^{j}\mathbf{t}_0||$ for neighboring camera pairs $(i, j)$, where ${}_{i}^{j}\mathbf{t}$ represents the translation from the camera $i$ to $j$.

## V. EXPERIMENTAL RESULTS

In order to evaluate the proposed system, we conduct extensive experiments with synthetic datasets along with real-world datasets. Using the synthetic datasets with ground-truth, we quantitatively measure the accuracy by computing the Root Mean Squared Error(RMSE) between the estimated poses and the ground truth. In addition, we compare average inlier ratio and average reprojection error to observe the tendency of experimental results. In the real datasets, we qualitatively evaluate the performance by overlaying the estimated trajectory to the satellite images. Additionally, we show the effectiveness of online extrinsic calibration in both synthetic and real datasets with comparative experiment.

| Sequences | Translation RMSE(m) | | | Average Inlier Ratio(%) | | | Average reprojection Error(°) | | | # of frame | Length(m) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | NoisyExt | OnlineExt | GTExt | NoisyExt | OnlineExt | GTExt | NoisyExt | OnlineExt | GTExt | | |
| 4view-MultiFOV | 30.01 | 0.60 | 0.23 | 28.1 | 47.8 | 51.0 | 0.60 | 0.19 | 0.13 | 2200 | 440 |
| Static Urban | 83.17 | 5.11 | 5.0 | 21.6 | 50.2 | 53.2 | 0.53 | 0.37 | 0.25 | 3000 | 2000 |
| Dynamic Urban | 65.65 | 5.68 | 1.62 | 27.9 | 56.9 | 57.5 | 0.48 | 0.29 | 0.18 | 1000 | 1320 |
| Cloudy Urban | 42.28 | 1.43 | 0.42 | 31.1 | 59.1 | 60.7 | 0.51 | 0.31 | 0.22 | 300 | 350 |
| Sunset Urban | 25.86 | 1.92 | 0.37 | 36.4 | 51.8 | 59.8 | 0.59 | 0.23 | 0.21 | 300 | 350 |

TABLE I: Quanitative evaluations with synthetic datasets. NoisyExt yields largely incorrect trajectories while OnlineExt shows good performance close to GTExt. See Section V-C for more details.
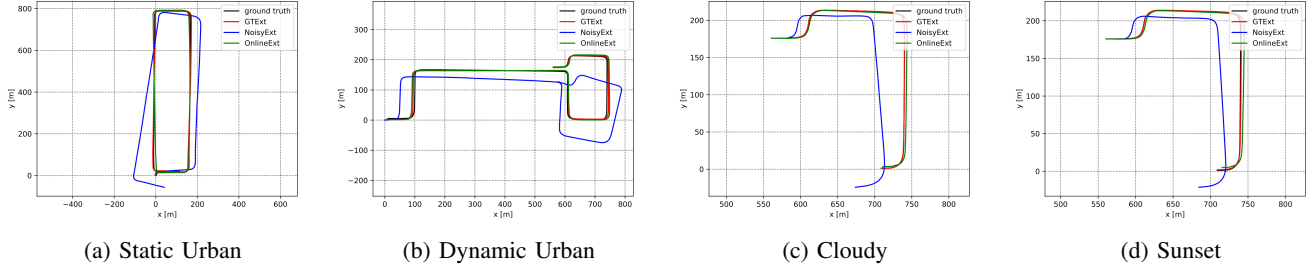


(a) Static Urban　　　　(b) Dynamic Urban　　　　(c) Cloudy　　　　(d) Sunset

Fig. 5: The estimated trajectories with the synthetic data-sets. The red, blue, green, and black lines represent the GTExt, OnlineExt, NoisyExt, and GT respectively. See Section V-C for more details.


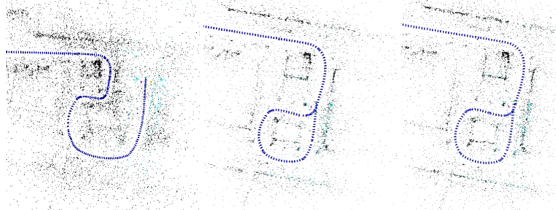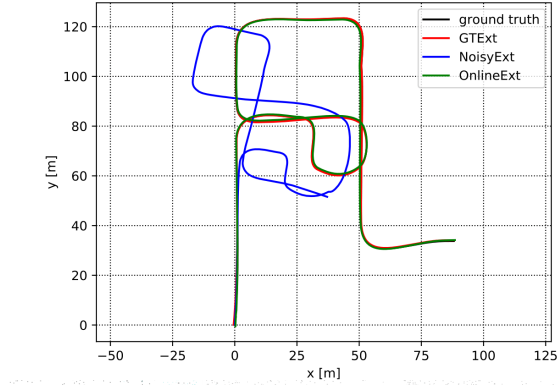
Fig. 6: Experimental results in 4view-MultiFOV sequence. Top: the estimated trajectories. Bottom: NoisyExt (left), OnlineExt (center), GTExt (right).

| | Original fisheye | Hybrid projection (ours) |
|---|---|---|
| Descriptor distance | 124.85 | 35.27 |
| Inlier ratio(%) | 21.53 | 73.16 |

TABLE II: Experiments of the hybrid projection model.

$3\times$, mainly by removing the lens distortion and aligning the projections. The boosted matching performance significantly contributes the robustness of VO systems.

### B. Synthetic and Real Datasets

We render four urban sequences with different structures lighting conditions using Blender. As a baseline we use 4view-MultiFOV which is modified from the urban canyon dataset by Zhang et al. [23]. Static Urban is a 2km-long sequence with a significant illumination change due to building shadows. Dynamic Urban is a 1.3km-long sequence with moving vehicles. Cloudy Urban and Sunset Urban are 350m long sequences of the same scene with different weather conditions. All images are rendered for four simulated fisheye cameras (220° FOV) with $1600 \times 1532$ resolution, which is same FOV and resolution with the real camera. The ground truth poses, camera intrinsic parameters and extrinsic parameters are included in the datasets.

For real datasets we use the four global-shutter camera rig on the vehicle as shown in Figure 1. All cameras output four software-synchronized $1600 \times 1532$ images at 10Hz. We use a standard camera rig calibration with a large checkerboard. The datasets are collected by driving the vicinity of Hanyang university, and they contain many challenges of harsh illumination changes, highly dynamic road with many moving vehicles, and narrow streets.

### C. Experiments with Synthetic Datasets

To evaluate the robustness and accuracy, we conduct an experiment by providing a randomly perturbed camera extrinsics to the system (zero-mean Gaussian noise with $\sigma = 5°$). NoisyExt is the result with the noisy extrinsics,

### A. Evaluation of Hybrid Projection Model

To test our hybrid projection model, we conduct two simple experiments. First, we compare the ORB descriptor similarity during intra-view tracking. The mean hamming distances of the 2000 features for 100 frames in the original fisheye images and in our hybrid-warped images are compared. Second, we compare the inlier ratios of feature correspondences between inter-view matching. The inlier ratio is computed only from the feature matches that satisfies the epipolar constraints.

As shown in Table II our projection model reduces the average descriptor distance and boosts the inlier ratio more than

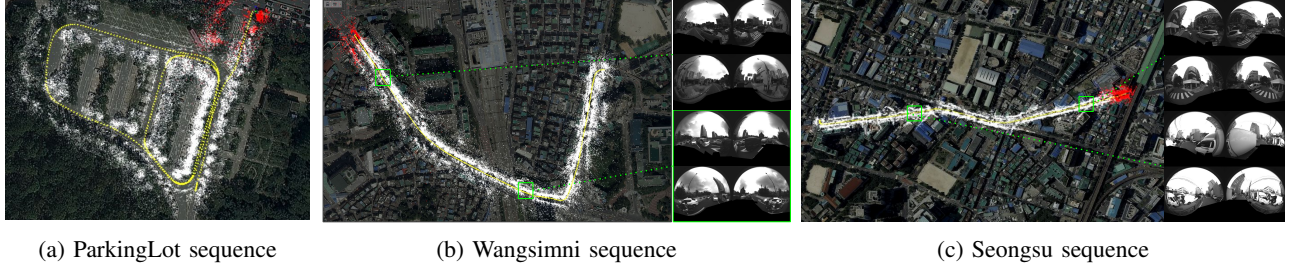(a) ParkingLot sequence      (b) Wangsimni sequence      (c) Seongsu sequence

Fig. 7: The estimated trajectories (yellow) of the real sequence overlaid on the satellite map. The white points represent the reconstructed 3D landmarks. See Section V-D for detailed discussion.
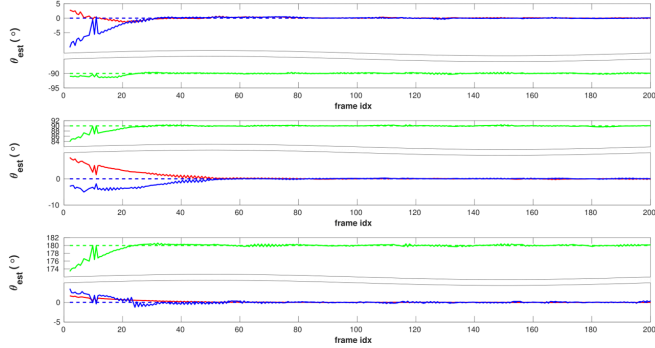


Fig. 8: The three plots correspond to each camera's relative pose to the first camera. Red, green, and blue lines represent pitch, roll, and yaw values respectively. The dotted line is the ground-truth. With online extrinsic calibration the initial noisy values quickly converges to the ground-truth.

and OnlineExt is the result with online extrinsic calibration. For comparison GT refers to the ground-truth rig trajectory, and GTExt is the VO result with the ground-truth extrinsics.

Quantitative and qualitative comparison is shown in Table I, Figure 5 and Figure 6. While the VO with noisy extrinsics fails to estimate correct trajectories, with online calibration overall error decreases drastically and the trajectory is also accurately estimated. We also observe that average inlier ratio and average reprojection error are significantly improved close to GTExt. Figure 8 gives in-depth view of online extrinsic calibration. Starting with noisy extrinsics, the extrinsic parameters are updated quickly and converges to the ground-truth within 100 frames.

### D. Experiments with Real Datasets

Among the collected datasets, we present the results of ParkingLot, Wangsimni, and Seongsu. ParkingLot has a loop trajectory and Figure 7(a) shows the accuracy of our system (note that our system do not use loop closing). Figure 7(b), (c) shows the Wangsimni and Seongsu results. Wangsimni sequence is taken in a heavy traffic with many moving objects. Our system is able to reject those outliers and robustly estimate the accurate trajectory. In Seongsu sequence the vehicle is driven in narrow passages, and the distance to the near buildings is very close compared to the camera baseline. The hybrid warping can successfully find the matching features in this challenging situations.
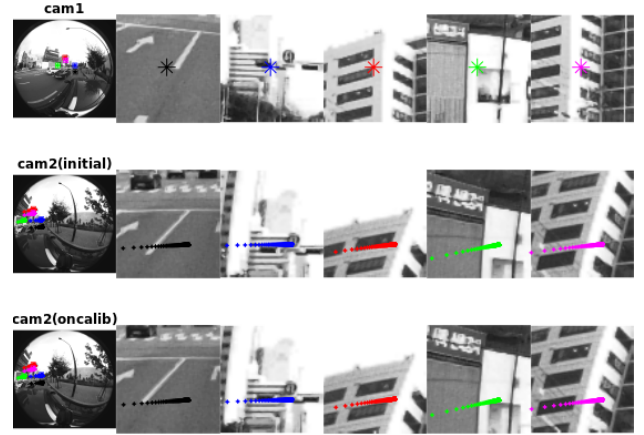


Fig. 9: Online extrinsic calibration result in a real dataset. Top: sample point locations in the reference camera. Middle: epipolar lines by initial extrinsics in the target camera. Bottom: corrected epipolar lines by online extrinsic calibration.

The online calibration makes a big contribution in real sequences. Figure 9 shows the reprojected epipolar lines of a feature in a neighbor view. With the initial calibration, the epipolar lines are quite off from the true matches, but with online calibration they fall on the correct position.

## VI. CONCLUSION

In this paper we propose a novel omnidirectional visual odometry system for a wide-baseline camera rig with wide-FOV lenses. To deal with the challenges from the fisheye distortion and appearance changes due to wide-baseline, we add a hybrid projection model, a multi-view P3P RANSAC algorithm, and online extrinsic calibration in local bundle adjustment. The extensive experimental evaluation using both synthetic datasets with ground-truth and real sequences verifies that the proposed components are effective in solving the problems.

# REFERENCES

[1] 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.

[2] R. Wang, M. Schwörer, and D. Cremers, "Stereo dso: Large-scale direct sparse visual odometry with stereo cameras."

[3] L. Kneip, H. Li, and Y. Seo, "Upnp: An optimal o(n) solution to the absolute pose problem with universal applicability," in *ECCV*, 2014.

[4] G. Schweighofer and A. Pinz, "Globally optimal o(n) solution to the pnp problem for general camera models," in *BMVC*, 2008.

[5] L. Kneip, D. Scaramuzza, and R. Siegwart, "A novel parametrization of the perspective-three-point problem for a direct computation of absolute camera position and orientation," 2011.

[6] 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.

[7] 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.

[8] C. Forster, M. Pizzoli, and D. Scaramuzza, "Svo: Fast semi-direct monocular visual odometry," in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 15–22.

[9] D. Caruso, J. Engel, and D. Cremers, "Large-scale direct slam for omnidirectional cameras." in *IROS*, vol. 1, 2015, p. 2.

[10] P. Liu, L. Heng, T. Sattler, A. Geiger, and M. Pollefeys, "Direct visual odometry for a fisheye-stereo camera," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 1746–1752.

[11] H. Matsuki, L. von Stumberg, V. Usenko, J. Stückler, and D. Cremers, "Omnidirectional dso: Direct sparse odometry with fisheye cameras," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3693–3700, 2018.

[12] G. Hee Lee, F. Faundorfer, and M. Pollefeys, "Motion estimation for self-driving cars with a generalized camera," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 2746–2753.

[13] L. Heng, G. H. Lee, and M. Pollefeys, "Self-calibration and visual slam with a multi-camera system on a micro aerial vehicle," *Autonomous Robots*, vol. 39, no. 3, pp. 259–277, 2015.

[14] P. Liu, M. Geppert, L. Heng, T. Sattler, A. Geiger, and M. Pollefeys, "Towards robust visual odometry with a multi-camera system."

[15] H. Guan and W. A. Smith, "Brisks: Binary features for spherical images on a geodesic grid," in *Proc. CVPR*, 2017, pp. 4886–4894.

[16] Q. Zhao, W. Feng, L. Wan, and J. Zhang, "Sphorb: A fast and robust binary feature on the sphere," *International Journal of Computer Vision*, vol. 113, no. 2, pp. 143–159, 2015.

[17] C. Häne, L. Heng, G. H. Lee, A. Sizov, and M. Pollefeys, "Real-time direct dense matching on fisheye images using plane-sweeping stereo," in *3D Vision (3DV), 2014 2nd International Conference on*, vol. 1. IEEE, 2014, pp. 57–64.

[18] W. Gao and S. Shen, "Dual-fisheye omnidirectional stereo," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 6715–6722.

[19] S. Agarwal, K. Mierle, and Others, "Ceres solver," http://ceres-solver.org.

[20] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, "Orb: An efficient alternative to sift or surf," in *Computer Vision (ICCV), 2011 IEEE international conference on*. IEEE, 2011, pp. 2564–2571.

[21] B. D. Lucas, T. Kanade, *et al.*, "An iterative image registration technique with an application to stereo vision," 1981.

[22] 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.

[23] Z. Zhang, H. Rebecq, C. Forster, and D. Scaramuzza, "Benefit of large field-of-view cameras for visual odometry," in *Robotics and Automation (ICRA), 2016 IEEE International Conference on*. IEEE, 2016, pp. 801–808.