

Supplementary Material for Paper 1640

Anonymous Author(s)

Submission Id: 1640

1 PRELIMINARIES OF THE SURROUND-VIEW SYSTEM

This section describes how to generate a surround-view from images captured by the cameras in the surround-view system.

Given the ground coordinate system O_G and a surround-view system consisting of multiple cameras, the pose of camera C_i is denoted by T_{C_iG} . For a point $P_G = [X_G, Y_G, Z_G, 1]^T$ in O_G , its corresponding pixel coordinate p_{C_i} in the camera coordinate system of C_i is given by,

$$p_{C_i} = \frac{1}{Z_{C_i}} K_{C_i} T_{C_iG} P_G \quad (1)$$

where Z_{C_i} is the depth of P_G in C_i 's coordinate system, and K_{C_i} is the intrinsic matrix of camera C_i , which can be estimated by Zhang's salient work [3] and some subsequent work of others [1, 4]. It's worth mentioning that the poses of cameras in the surround-view system are usually determined by offline calibration. In our solution, the scheme proposed by Shao *et al.* in [2] is adopted.

Consider a point $p_G = [u_G, v_G, 1]^T$ on the bird's-eye-view image. Its corresponding point on the ground plane is denoted by $P_G = [X_G, Y_G, Z_G = 0]^T$ with respect to the ground coordinate system. The relationship between p_G and P_G can be represented as,

$$p_G = K_G P_G \quad (2)$$

and the transformation matrix K_G is defined as,

$$K_G = \begin{bmatrix} \frac{1}{d_{X_G}} & 0 & \frac{W}{2d_{X_G}} \\ 0 & -\frac{1}{d_{Y_G}} & \frac{H}{2d_{Y_G}} \\ 0 & 0 & 1 \end{bmatrix} \quad (3)$$

where d_{X_G} and d_{Y_G} define the physical size of each pixel¹, and W and H are the width and height of the synthesized surround-view, respectively. It is worth mentioning that since $Z_G = 0$, it is ignored implicitly here. By combining Eq. 1 and Eq. 2, we can get,

$$p_{C_i} = \frac{1}{Z_{C_i}} K_{C_i} T_{C_iG} K_G^{-1} p_G \quad (4)$$

Eq. 4 actually depicts the relationship of a point p_{C_i} on the image plane of camera C_i and its projection p_G on the surround-view. From Eq. 4, we can generate a bird's-eye-view image by projecting the undistorted image of camera C_i onto the ground,

$$I_{GC_i}(p_G) = I_{C_i}(p_{C_i}) \quad (5)$$

where I_{C_i} is the undistorted fisheye image captured by camera C_i , and I_{GC_i} is the ground projection of I_{C_i} , namely the bird's-eye-view image. Then, the surround-view image can be synthesized with appropriate stitching seams.

¹More accurately, each pixel in the surround-view image corresponds to a $d_{X_G} \times d_{Y_G}$ physical area on the ground plane.

2 JACOBIANS OF THE BI-CAMERA ERROR

As mentioned in the manuscript, the bi-camera error term $\epsilon_{p_G}^{bi}$ of a point p_G on the surround-view is defined as,

$$\epsilon_{p_G}^{bi} = \frac{1}{|\mathcal{P}|} \sum_{p_s \in \mathcal{P}} I_{C_i} \left(\lambda_{p_G}^{C_i} K_{C_i} \exp \left(\xi_{C_iG}^\wedge \right) K_G^{-1} p_G + p_s \right) - \gamma_{ij} I_{C_j} \left(\lambda_{p_G}^{C_j} K_{C_j} \exp \left(\xi_{C_jG}^\wedge \right) K_G^{-1} p_G \right) \quad (6)$$

where I_{C_i} and I_{C_j} are undistorted images captured by C_i and C_j , respectively. K_{C_i} and K_{C_j} are intrinsics of C_i and C_j , respectively. ξ_{C_iG} and ξ_{C_jG} are poses of C_i and C_j in Lie algebra form, respectively. K_G stands for the transformation matrix from the surround-view coordinate system to the ground coordinate system. \mathcal{P} is a set that contains the relative pixel coordinates of all the utilized points to p_{C_i} , and is defined as,

$$\mathcal{P} = \{[i, j]^T \mid i, j = -2, 0, 2\}. \quad (7)$$

Then, in this section, Jacobians of the bi-camera error term to both the camera pose and the inverse depth will be deduced in detail.

Jacobian to the pose. The Jacobian J_p of the bi-camera error term $\epsilon_{p_G}^{bi}$ to camera C_i 's pose ξ_{C_iG} can be expressed as,

$$J_p = \frac{\partial \epsilon_{p_G}^{bi}}{\partial \xi_{GC_i}^T}. \quad (8)$$

It can be decomposed to 4 parts with the chain rule,

$$J_p = \frac{\partial \epsilon_{p_G}^{bi}}{\partial I_{C_i}} \cdot \frac{\partial I_{C_i}}{\partial p_{C_i}^T} \cdot \frac{\partial p_{C_i}}{\partial p_{C_i}^T} \cdot \frac{\partial p_{C_i}}{\partial \xi_{C_iG}^T}. \quad (9)$$

Next, we discuss these 4 parts one by one.

(1) $\partial \epsilon_{p_G}^{bi} / \partial I_{C_i}$ is the derivative of the error $\epsilon_{p_G}^{bi}$ to pixel intensities of image I_{C_i} . Actually, from Eq. 6, it's easy to know that this term is equal to one,

$$\frac{\partial \epsilon_{p_G}^{bi}}{\partial I_{C_i}} = 1. \quad (10)$$

(2) $\partial I_{C_i} / \partial p_{C_i}^T$ is the average intensity gradient, which is generally computed by the Sobel operator, of image I_{C_i} at all the pixels in the local window \mathcal{P} whose center is p_{C_i} . Actually, this term can also be approximated just by the intensity gradient at p_{C_i} (the window of the Sobel operator needs to be enlarged accordingly). Thus, $\partial I_{C_i} / \partial p_{C_i}^T$ can be given as,

$$\frac{\partial I_{C_i}}{\partial p_{C_i}^T} = \begin{bmatrix} \frac{\partial I_{C_i}}{\partial u_{C_i}} & \frac{\partial I_{C_i}}{\partial v_{C_i}} \end{bmatrix} \triangleq \begin{bmatrix} \nabla I_{C_i}^{u_{C_i}} & \nabla I_{C_i}^{v_{C_i}} \end{bmatrix} \quad (11)$$

where u_{C_i} and v_{C_i} are both coordinate values of p_{C_i} .

(3) $\partial p_{C_i} / \partial p_{C_i}^T$ is the derivative of a pixel's 2D coordinate to its 3D position in the camera coordinate system. From the pin-hole

camera model, we have

$$\frac{\partial \mathbf{P}_{C_i}}{\partial \mathbf{P}_{C_i}^T} = \begin{bmatrix} \frac{f_x^i}{Z_{C_i}} & 0 & -\frac{f_x^i X_{C_i}}{Z_{C_i}^2} \\ 0 & \frac{f_y^i}{Z_{C_i}} & -\frac{f_y^i Y_{C_i}}{Z_{C_i}^2} \end{bmatrix} \quad (12)$$

where f_x^i and f_y^i are focal lengths of C_i , and X_{C_i} , Y_{C_i} and Z_{C_i} are coordinate values in three axes of \mathbf{P}_{C_i} in C_i 's coordinate system.

(4) $\partial \mathbf{P}_{C_i} / \partial \xi_{C_i G}^T$ is the derivative of the 3D point \mathbf{P}_{C_i} to the camera pose $\xi_{C_i G}$,

$$\frac{\partial \mathbf{P}_{C_i}}{\partial \xi_{C_i G}^T} = \begin{bmatrix} \mathbf{I}_{3 \times 3} & -\mathbf{P}_{C_i}^\wedge \end{bmatrix} \quad (13)$$

where \mathbf{I} is a 3×3 identity matrix and $\mathbf{P}_{C_i}^\wedge$ is the 3×3 anti-symmetric matrix generated from \mathbf{P}_{C_i} . By merging the four terms in Eqs. 10~13, we can get the final form of the Jacobian \mathbf{J}_p ,

$$\mathbf{J}_p = \begin{bmatrix} \nabla \mathbf{I}_{C_i}^{u_{C_i}} & \nabla \mathbf{I}_{C_i}^{v_{C_i}} \end{bmatrix} \begin{bmatrix} \frac{f_x^i}{Z_{C_i}} & 0 & -\frac{f_x^i X_{C_i}}{Z_{C_i}^2} \\ 0 & \frac{f_y^i}{Z_{C_i}} & -\frac{f_y^i Y_{C_i}}{Z_{C_i}^2} \end{bmatrix} \begin{bmatrix} \mathbf{I}_{3 \times 3} & -\mathbf{P}_{C_i}^\wedge \end{bmatrix}. \quad (14)$$

Jacobian to the inverse depth. The Jacobian \mathbf{J}_d of the bi-camera error term $\varepsilon_{p_G}^{bi}$ to point \mathbf{p}_{C_j} 's inverse depth $\lambda_{p_G}^{C_j}$ can be expressed as,

$$\mathbf{J}_d = \frac{\partial \varepsilon_{p_G}^{bi}}{\partial \lambda_{p_G}^{C_j}}. \quad (15)$$

With the chain rule, it can also be decomposed as,

$$\mathbf{J}_d = \frac{\partial \varepsilon_{p_G}^{bi}}{\partial \mathbf{P}_{C_i}^T} \cdot \frac{\partial \mathbf{P}_{C_i}}{\partial \mathbf{P}_{C_j}^T} \cdot \frac{\partial \mathbf{P}_{C_j}}{\partial \lambda_{p_G}^{C_j}}. \quad (16)$$

Next, these three simpler parts are discussed one by one.

(1) $\partial \varepsilon_{p_G}^{bi} / \partial \mathbf{P}_{C_i}^T$ is the derivative of the error $\varepsilon_{p_G}^{bi}$ to \mathbf{p}_G 's corresponding 3D position \mathbf{P}_{C_i} in C_i 's camera coordinate system. This term can be obtained by combining Eqs. 10 ~ 12, which is given as,

$$\frac{\partial \varepsilon_{p_G}^{bi}}{\partial \mathbf{P}_{C_i}^T} = \begin{bmatrix} \nabla \mathbf{I}_{C_i}^{u_{C_i}} & \nabla \mathbf{I}_{C_i}^{v_{C_i}} \end{bmatrix} \begin{bmatrix} \frac{f_x^i}{Z_{C_i}} & 0 & -\frac{f_x^i X_{C_i}}{Z_{C_i}^2} \\ 0 & \frac{f_y^i}{Z_{C_i}} & -\frac{f_y^i Y_{C_i}}{Z_{C_i}^2} \end{bmatrix}. \quad (17)$$

(2) $\partial \mathbf{P}_{C_i} / \partial \mathbf{P}_{C_j}^T$ is the derivative of \mathbf{P}_{C_i} 's 3D coordinate in C_i 's coordinate system to its corresponding 3D point in C_j 's coordinate system. This term is given as,

$$\frac{\partial \mathbf{P}_{C_i}}{\partial \mathbf{P}_{C_j}^T} = \mathbf{T}_{C_i G} \mathbf{T}_{C_j G}^{-1} = \mathbf{T}_{C_i C_j}. \quad (18)$$

where $\mathbf{T}_{C_i C_j}$ is the relative pose of C_j to C_i .

(3) $\partial \mathbf{P}_{C_j} / \partial \lambda_{p_G}^{C_j}$ is the derivative of a 3D point \mathbf{P}_{C_j} to its inverse depth. It can be expressed as,

$$\frac{\partial \mathbf{P}_{C_j}}{\partial \lambda_{p_G}^{C_j}} = -\frac{1}{(\lambda_{p_G}^{C_j})^2} \mathbf{K}_{C_j}^{-1} \mathbf{p}_{C_j} = -\frac{1}{(\lambda_{p_G}^{C_j})} \mathbf{P}_{C_j}. \quad (19)$$

Thus, the final form of the Jacobian \mathbf{J}_d is given by,

$$\begin{aligned} \mathbf{J}_d &= -\frac{1}{(\lambda_{p_G}^{C_j})} \begin{bmatrix} \nabla \mathbf{I}_{C_i}^{u_{C_i}} & \nabla \mathbf{I}_{C_i}^{v_{C_i}} \end{bmatrix} \begin{bmatrix} \frac{f_x^i}{Z_{C_i}} & 0 & -\frac{f_x^i X_{C_i}}{Z_{C_i}^2} \\ 0 & \frac{f_y^i}{Z_{C_i}} & -\frac{f_y^i Y_{C_i}}{Z_{C_i}^2} \end{bmatrix} \mathbf{T}_{C_i C_j} \mathbf{P}_{C_j} \\ &= -\frac{1}{(\lambda_{p_G}^{C_j})} \begin{bmatrix} \nabla \mathbf{I}_{C_i}^{u_{C_i}} & \nabla \mathbf{I}_{C_i}^{v_{C_i}} \end{bmatrix} \begin{bmatrix} \frac{f_x^i}{Z_{C_i}} & 0 & -\frac{f_x^i X_{C_i}}{Z_{C_i}^2} \\ 0 & \frac{f_y^i}{Z_{C_i}} & -\frac{f_y^i Y_{C_i}}{Z_{C_i}^2} \end{bmatrix} \mathbf{P}_{C_i}. \end{aligned} \quad (20)$$

As all derivative relationships between the bi-camera error term and the optimized variables have been deduced, the objective function of ROECS can then be minimized with any non-linear optimization scheme and thereby we can get accurate extrinsics of the SVS.

REFERENCES

- [1] Fenglei Du and Michael Brady. 1993. Self-calibration of the Intrinsic Parameters of Cameras for Active Vision Systems. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR'93)*. IEEE, New York, USA, 477–482. <https://doi.org/10.1109/CVPR.1993.341087>
- [2] Xuan Shao, Xiao Liu, Lin Zhang, Shengjie Zhao, Ying Shen, and Yukai Yang. 2019. Revisit Surround-view Camera System Calibration. In *International Conference on Multimedia and Expo (ICME'19)*. IEEE, Shanghai, China, 1486–1491. <https://doi.org/10.1109/ICME.2019.00257>
- [3] Zhengyou Zhang. 1999. Flexible Camera Calibration by Viewing a Plane from Unknown Orientations. In *IEEE International Conference on Computer Vision (ICCV'99)*. IEEE, Kerkyra, Greece, 666–673. <https://doi.org/10.1109/ICCV.1999.791289>
- [4] Haijiang Zhu, Jinfu Yang, and Zhongtian Liu. 2009. Fisheye Camera Calibration with Two Pairs of Vanishing Points. In *International Conference on Information Technology and Computer Science (ITCS'09)*. IEEE, Kiev, Ukraine, 321–324. <https://doi.org/10.1109/ITCS.2009.72>