Radially-Distorted Conjugate Translations Supplementary Material

James Pritts¹ Zuzana Kukelova¹

Viktor Larsson² Ondřej Chum¹

Visual Recognition Group, CTU in Prague¹

Centre for Mathematical Sciences, Lund University²

A. Transfer Error

The scene plane is tessellated by a 10x10 square grid of points, denoted $\{X_i\}$, with a 1 meter spacing between adjacent points. Suppose that $\mathbf{y} \leftrightarrow \mathbf{y}'$ is an undistorted point correspondence induced by the conjugate translation $H_{\mathbf{u}} = [I_3 + \mathbf{u} \mathbf{l}^{\top}]$ in the imaged scene plane (here we assume that $s_i^{\mathbf{u}} = 1$ since we speak about an individual point correspondence).

Points { X_i } are translated by 1 meter on the scene plane in the direction of translation induced by the preimage of the point correspondence $\mathbf{y} \leftrightarrow \mathbf{y}'$ giving the translated grid $\{\mathbf{X}'_i\}$. The purpose of constructing the grid and it's translation is to uniformly cover the scene plane that the camera images in its field of view. In this way, the accuracy of the conjugate translation and lens-distortion parameter estimation can be measured across most of the image. The conjugate translation H_u is not used directly because the magnitude of translation may span the extent of the scene plane, so applying it to the tessellation would transform the grid out of the field of view.

Let the camera be parameterized by the camera matrix $P = (AH)^{-1}$ (see Sec. 5.4 for the definition of the camera matrix) that pointwise maps the scene plane Π to the imaged scene plane π and division model parameter λ . The preimages of the undistorted point correspondence $\mathbf{y} \leftrightarrow \mathbf{y}'$ in the scene-plane coordinate system is, respectively, $\beta \mathbf{Y} = \mathbf{P}^{-1}\mathbf{y}$ and $\beta' \mathbf{Y}' = \mathbf{P}^{-1}\mathbf{y}'$. The translation t of the preimages in the scene plane coordinate system is $\mathbf{t} = \mathbf{Y} - \mathbf{Y}' = \left(t_x, t_y, 0\right)^\top.$

Then $\|\mathbf{t}\|$ is the magnitude of translation between the repeated scene elements in the scene-plane coordinate system. Denote the homogeneous translation matrix T(t) to be the matrix constructed from t as

$$\mathbf{T}(\mathbf{t}) = \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix}.$$
 (12)

The translation of the grid points by unit distance is given by $\mathbf{X}'_i = T(\mathbf{t}/||\mathbf{t}||)\mathbf{X}_i$. Recall from (1) that a conjugate translation has the form $PT(\cdot)P^{-1}$. Using (2), the conjugate translation of unit distance in the direction of point correspondences $\mathbf{y} \leftrightarrow \mathbf{y}'$ is

$$\begin{aligned} \mathbf{H}_{\mathbf{u}/\|\mathbf{t}\|} &= \mathbf{P}\mathbf{I}_{3}\mathbf{P}^{-1} + \mathbf{P}\begin{pmatrix} t_{x}/\|\mathbf{t}\|\\ t_{y}/\|\mathbf{t}\|\\ 1 \end{pmatrix} \begin{bmatrix} \mathbf{P}^{-\top} \begin{pmatrix} 0\\0\\1 \end{pmatrix} \end{bmatrix}^{\top} \\ &= [\mathbf{I}_{3} + \frac{\mathbf{u}}{\|\mathbf{t}\|}\mathbf{I}^{\top}]. \end{aligned}$$
(13)

The unit conjugate translation $H_{\mathbf{u}/\|\mathbf{t}\|}$ can be written in terms of the conjugate translation H_u induced by the undistorted point correspondence $\mathbf{y} \leftrightarrow \mathbf{y}'$ as

$$\mathbf{I}_{3} + \frac{\mathbf{u}}{\|\mathbf{t}\|} \mathbf{l}^{\mathsf{T}} = \mathbf{I}_{3} + \frac{1}{\|\mathbf{t}\|} [\mathbf{I}_{3} + \mathbf{u} \mathbf{l}^{\mathsf{T}} - \mathbf{I}_{3}]$$
$$= \mathbf{I}_{3} + \frac{1}{\|\mathbf{t}\|} [\mathbf{H}_{\mathbf{u}} - \mathbf{I}_{3}].$$
(14)

The derivation of (14) gives the form of transformation used in the transfer error $\Delta_{\rm RMS}^{\rm xfer}$ defined in Sec. 5.1, which maps from the undistorted points of the grid $\{\mathbf{x}_i\}$ to their translated correspondences $\{\mathbf{x}'_i\}$.

B. Computational Complexity.

Table B.1 lists the elimiation template sizes for the proposed solvers. The average time to compute the solutions for a given input for a solver is directly proportional to the elimination template size. The solvers are implemented in MATLAB and C++. Significantly faster implementations are possible with a pure C++ port. Still the solvers are sufficiently fast for use in RANSAC. The proposed solvers have an average solve time from 0.3 to 2 milliseconds.

$\mathrm{H}2.5\mathrm{lu}\lambda$	$\mathrm{H3lu} s_{\mathbf{u}} \lambda$	$\mathrm{H3.5}\mathbf{luv}\lambda$	$\mathtt{H4}\mathbf{luv}s_{\mathbf{v}}\lambda$
14x18	24x26	54x60	76x80

Table B.1: Template sizes for the proposed solvers.

C. Extended Experiments

The extended real-data experiments in the following pages include (i) images with lesser radial distortion from consumer cameras and mobile phones; the images also demonstrate the proposed method's effectiveness on diverse scene content, (ii) images for very wide field-of-view lenses (8mm and 12mm), (iii) and an additional local-optimization experiment similar to the one in Fig. 6 on a GoPro Hero 4 image taken with its medium field-of-view setting, which further demonstrates the need for a minimal solver that jointly estimates lens distortion with affine-rectification to achieve an accurate undistortions and rectifications.



Figure C.1: Narrow field-of-view and diverse scene-content experiments for H2.5lu λ +LO. The proposed method works well if the input image has little or no radial lens distortion. This imagery is typical of consumer cameras and mobile phone cameras. The images are diverse and contain unconventional scene content. Input images are on the top row; undistorted images are on the middles row, and the rectified images are on the bottom row.



Figure C.2: GoPro Hero 4 at the medium setting for different solvers. Results from LO-RANSAC (see Sec. 4) for H2lu, which omits distortion, and the proposed solvers H2.5lu λ and H3.5luv λ . The top row has rectifications after local optimization (LO); The bottom row has undistortions estimated from the best *minimal* sample. LO-RANSAC fails from the poor initializations by H2lu (column 2). The proposed solvers in columns 3 and 4 give a correct rectification. The bottom left has a chessboard undistorted using the division parameter estimated from the building facade by H2.5lu λ +LO.



Figure C.3: Very wide-angle images undistorted and rectified with $H2.5lu\lambda+LO$. The left column is an image from an 8mm lens, and the right column is from a 12mm lens. The top row contains the input images; the middle row contains the undistorted images, and the bottom row contains the rectified images. The division model [7] used for radial lens distortion has only 1 parameter, which may impose limits for modeling extreme lens distortion.



Figure C.4: *Problem difficulty and method robustness* The input to the method is ungrouped affine-covriant features (top left). Common problems include missed detections of repeated texture, duplicate detections, and detections due to compression artifacts, all of which are visible in this example. The inliers with respect to the total number of detected features can be a very small proportion (top right). Still the method can estimate accurate undistortion (bottom left) and affine rectification (bottom right) even from a very sparse sampling of the inlying affine-covariant features on the scene plane (top right).