# Rethinking Generic Camera Models for Deep Single Image Camera Calibration to Recover Rotation and Fisheye Distortion
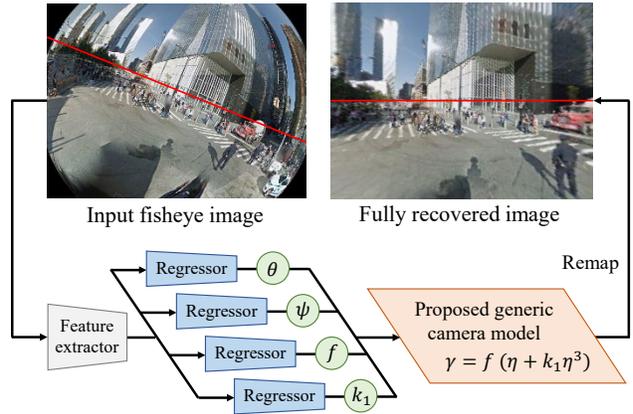
Nobuhiko Wakai[1]    Satoshi Sato[1]    Yasunori Ishii[1]    Takayoshi Yamashita[2]

[1] Panasonic Corporation    [2] Chubu University

{*lastname.firstname*}@jp.panasonic.com    takayoshi@isc.chubu.ac.jp

## Abstract

*Although recent learning-based calibration methods can predict extrinsic and intrinsic camera parameters from a single image, the accuracy of these methods is degraded in fisheye images. This degradation is caused by mismatching between the actual projection and expected projection. To address this problem, we propose a generic camera model that has the potential to address various types of distortion. Our generic camera model is utilized for learning-based methods through a closed-form numerical calculation of the camera projection. Simultaneously to recover rotation and fisheye distortion, we propose a learning-based calibration method that uses the camera model. Furthermore, we propose a loss function that alleviates the bias of the magnitude of errors for four extrinsic and intrinsic camera parameters. Extensive experiments demonstrated that our proposed method outperformed conventional methods on two large-scale datasets and images captured by off-the-shelf fisheye cameras. Moreover, we are the first researchers to analyze the performance of learning-based methods using various types of projection for off-the-shelf cameras.*

## 1. Introduction

Learning-based perception methods are widely used for surveillance, cars, drones, and robots. These methods are well established for many computer vision tasks. Most computer vision tasks require undistorted images; however, fisheye images have the superiority of a large field of view in visual surveillance [11], object detection [34], pose estimation [5], and semantic segmentation [26]. To use fisheye cameras through removing distortion, camera calibration is a desirable step before perception. Camera calibration is a long-studied topic in areas of computer vision, such as image undistortion [24, 47], image remapping [42], virtual object insertion [17], augmented reality [3], and stereo measurement [28]. In camera calibration, we cannot escape the trade-off between accuracy and usability that we need a cal-



**Figure 1.** Concept illustrations of our work. Our network predicts parameters in our proposed generic camera model to obtain fully recovered images using remapping. Red lines indicate horizontal lines in each of the images.

ibration object; hence, tackling the trade-off has been an open challenge, which we explain further in the following.

Calibration methods are classified into two categories: geometric-based and learning-based methods. Geometric-based calibration methods achieve high accuracy, but they require a calibration object, such as a cube [41] and planes [48], to obtain a strong geometric constraint. By contrast, learning-based methods can calibrate cameras without a calibration object from a general scene image [24, 47], which is called deep single image camera calibration. Although learning-based methods do not require a calibration object, the accuracy of these methods is degraded for fisheye images because of the mismatch between the actual projection and expected projection in conventional methods. In particular, calibration methods [29, 42] that predict both camera rotation and distortion have much room for improvement regarding addressing complex fisheye distortion. López-Antequera's method [29] was designed for non-fisheye cameras with radial distortion and cannot process fisheye distortion. Although four standard camera models are used for fisheye cameras, Wakai's method [42] supports

only one fisheye camera model.

Based on the observations above, we propose a new generic camera model for various fisheye cameras. The proposed generic camera model has the potential to address various types of distortion. For the generic camera model, we propose a learning-based calibration method that predicts extrinsic parameters (tilt and roll angles), focal length, and a distortion coefficient simultaneously from a single image, as shown in Figure 1. Our camera model is utilized for learning-based methods through a closed-form numerical calculation of camera projection. To improve the prediction accuracy, we use a joint loss function composed of each loss for the four camera parameters. Unlike heuristic approaches in conventional methods, our loss function makes significant progress; that is, we can determine the optimal joint weights based on the magnitude of errors for these camera parameters instead of the heuristic approaches.

To evaluate the proposed method, we conducted extensive experiments on two large-scale datasets [6, 30] and images captured by off-the-shelf fisheye cameras. This evaluation demonstrated that our method meaningfully outperformed conventional geometric-based [37] and learning-based methods [7, 24, 29, 42, 47]. The major contributions of our study are summarized as follows:

- We propose a learning-based calibration method for recovering camera rotation and fisheye distortion using the proposed generic camera model that has an adaptive ability for off-the-shelf fisheye cameras. To the best of our knowledge, we are the first researchers to calibrate extrinsic and intrinsic parameters of generic camera models from a single image.

- We propose a new loss function that alleviates the bias of the magnitude of errors between the ground-truth and predicted camera parameters for four extrinsic and intrinsic parameters to obtain accurate camera parameters.

- We first analyze the performance of learning-based methods using various off-the-shelf fisheye cameras. In previous studies, these conventional learning-based methods were evaluated using only synthetic images.

## 2. Related work

**Camera calibration:** Camera calibration estimates parameters composed of extrinsic parameters (rotation and translation) and intrinsic parameters (image sensor and distortion parameters). Geometric-based calibration methods have been developed using a strong constraint based on the calibration object [41, 48] or line detection [2, 37]. This constraint explicitly represents the relation between world coordinates and image coordinates for the stable optimization of calibration. By contrast, learning-based methods based

on convolutional neural networks calibrate cameras from a single image in the wild. In this study, we focus on learning-based calibration methods and describe them below.

Calibration methods for only extrinsic parameters have been proposed that are aimed at narrow view cameras [19, 32, 38, 39, 44, 45] and panoramic 360° images [10]. These methods cannot calibrate intrinsic parameters, that is, they cannot remove distortion. For extrinsic parameters and focal length, narrow-view camera calibration was developed with depth estimation [8, 15] and room layout [35]. These methods are not suitable for fisheye cameras because fisheye distortion is not negligible because the projection of the field of view is over 180°.

To address large distortion, calibration methods for only undistortion have been proposed that use specific image features, that is, segmentation information [47], straight lines [46], and ordinal distortion of part of the images [24]. Furthermore, Chao *et al*. [7] proposed undistortion networks based on generative adversarial networks [13]. These methods can process only undistortion and image remapping tasks.

For both extrinsic and intrinsic parameters, López-Antequera *et al*. [29] proposed a pioneering method for non-fisheye cameras. This method estimates distortion using a polynomial function model of perspective projection similar to Tsai's quartic polynomial model [41]. This polynomial function of the distance from a principal point has the two coefficients for second- and fourth-order terms. The method is only trainable for the second-order coefficient, and the fourth-order coefficient is calculated using a quadratic function of the second-order one. This method does not calibrate fisheye cameras effectively because the camera model does not represent fisheye camera projection. Additionally, Wakai *et al*. [42] proposed a calibration method for extrinsic parameters and focal length in fisheye cameras. Although four types of standard fisheye projection are used for camera models, for example, equisolid angle projection, the method of Wakai *et al*. [42] only expects equisolid angle projection. As discussed above, conventional learning-based calibration methods do not fully calibrate extrinsic and intrinsic parameters of generic camera models from a single image.

**Exploring loss landscapes:** To optimize networks effectively, loss landscapes have been explored after training [9, 14, 23] and during training [16]. In learning-based calibration methods, we have the problem that joint weights are difficult to determine without training. To stabilize training or the merging of heterogeneous loss components, the joint loss function was often defined [24, 29, 42, 47]. However, these joint weights were defined using experiments or the same values, that is, unweighted joints. These joint weights are hyperparameters that depend on networks and datasets. A hyperparameter search method was proposed by Akiba *et al*. [1]. However, hyperparameter search tools re-

quire large computational cost because they execute various conditions. Additionally, to analyze the optimizers, Goodfellow *et al.* [14] proposed an examination method for loss landscapes that use linear interpolation from the initial network weights to the final weights. To overcome the saddle points of loss landscapes, Dauphin *et al.* [9] proposed an optimization method based on Newton's method. Furthermore, Li *et al.* [23] developed a method for visualizing loss landscapes. Although these methods can explore high-order loss landscapes, the optimal values of joint loss weights have not been determined in learning-based calibration methods. Moreover, the aforementioned methods cannot explore loss landscapes without training because they require training results.

## 3. Proposed method

First, we describe our proposed camera model based on a closed-form solution for various fisheye cameras. Second, we describe our learning-based calibration method for recovering rotation and fisheye distortion. Finally, we explain a new loss function, with its notation and mechanism.

### 3.1. Generic camera model

Camera models are composed of extrinsic parameters $[\mathbf{R} \mid \mathbf{t}]$ and intrinsic parameters, and these camera models represent the mapping from world coordinates $\tilde{\mathbf{p}}$ to image coordinates $\tilde{\mathbf{u}}$ in homogeneous coordinates. This projection can be expressed for radial distortion [33] and fisheye models [40] as

$$\tilde{\mathbf{u}} = \begin{bmatrix} \gamma/d_u & 0 & c_u \\ 0 & \gamma/d_v & c_v \\ 0 & 0 & 1 \end{bmatrix} [\mathbf{R} \mid \mathbf{t}] \tilde{\mathbf{p}}, \qquad (1)$$

where $\gamma$ is distortion, $(d_u, d_v)$ is the image sensor pitch, $(c_u, c_v)$ is a principal point, $\mathbf{R}$ is a rotation matrix, and $\mathbf{t}$ is a translation vector. The subscripts of $u$ and $v$ denote the horizontal and vertical direction, respectively.

The generic camera model including fisheye lenses [12] is defined as

$$\gamma = \tilde{k_1}\eta + \tilde{k_2}\eta^3 + \cdots, \qquad (2)$$

where $\tilde{k_1}, \tilde{k_2}, \ldots$ are distortion coefficients, and $\eta$ is an incident angle. Note that the focal length is not defined explicitly, that is, the focal length is set to 1 mm, and the distortion coefficients represent distortion and implicit focal length.

### 3.2. Proposed camera model

A generic camera model with high order has the potential to achieve high calibration accuracy. However, this high-order function leads to unstable optimization, particularly for learning-based methods. Considering this problem, we propose a generic camera model for learning-based fisheye

**Table 1.** Comparison of absolute errors in fisheye camera models

| Reference model[1] | Mean absolute error [pixel] | | | |
| --- | --- | --- | --- | --- |
| | STG | EQD | ESA | ORT |
| Stereographic (STG) | – | 9.33 | 13.12 | 93.75 |
| Equidistance (EQD) | 9.33 | – | 3.79 | 23.58 |
| Equisolid angle (ESA) | 13.12 | 3.79 | – | 14.25 |
| Orthogonal (ORT) | 93.75 | 23.58 | 14.25 | – |
| Proposed generic model | **0.54** | **0.00** | **0.02** | **0.35** |

[1] Each reference model is compared with other fisheye models

calibration using the explicit focal length, given by

$$\gamma = f(\eta + k_1\eta^3), \qquad (3)$$

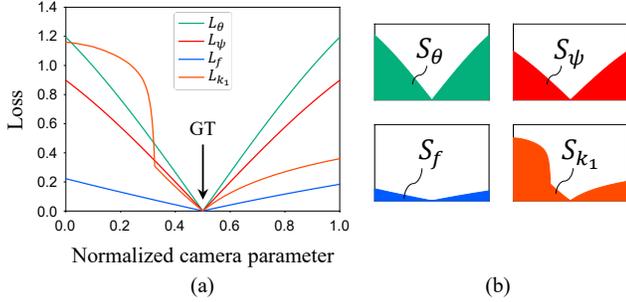where $f$ is the focal length and $k_1$ is a distortion coefficient.

**Evaluating our camera model:** Our generic camera model is a third-order polynomial function that corresponds to Taylor's expansion of the trigonometric function in fisheye cameras, that is, stereographic projection, equidistance projection, equisolid angle projection, and orthogonal projection. In the following, we show that our model can express trigonometric function models with slight errors.

We compared projection function, $\gamma = g(\eta)$, of the four trigonometric function models and our generic camera model, as shown in Table 1. In this comparison, we calculated mean absolute errors $\epsilon$ between pairs of the projection function $g_1$ and $g_2$. We defined the errors as $\epsilon = 1/(\pi/2) \int_0^{\pi/2} | g_1(\eta) - g_2(\eta) | \, d\eta$. This mean absolute errors simply represents mean distance errors in image coordinates. Our model is useful for various fisheye models because our model had the smallest mean absolute errors among the camera models in Table 1.

**Calculation easiness:** For our generic camera model, it is easy to calculate back-projection, which converts image coordinates to corresponding incident angles. When using back-projection, we must solve the generic camera model against incident angles $\eta$ in Equation (3). Practically, we can solve equations on the basis of iterations or closed-form. Non-fisheye cameras often use the iteration approaches [4]. By contrast, we cannot use the iteration approaches for fisheye cameras because large distortion prevents us from obtaining the initial values close to solutions. We therefore use a closed-form approach because the Abel-Ruffini theorem [49] shows that fourth-order or less algebraic equations are solvable.

### 3.3. Proposed calibration method

To calibrate various fisheye cameras, we propose a learning-based calibration method that uses our generic camera model. We use DenseNet-161 [18] pretrained on ImageNet [36] to extract image features and details as follows: First, we convert the image features using global average pooling [25] for regressors. Second, four individual regressors predict the normalized parameters (from 0 to 1) of

**Figure 2.** Difference between the non-grid bearing loss functions [42] for the camera parameters. (a) Each loss landscape along the normalized camera parameters using a predicted camera parameter with a subscript parameter and ground-truth parameters for the remaining parameters, and the ground-truth values are set to 0.5. (b) Areas $S$ integrating $L$ that shows (a) for the interval 0 to 1 for $\theta$, $\psi$, $f$, and $k_1$.

the tilt angle $\theta$, roll angle $\psi$, focal length $f$, and a distortion coefficient $k_1$. Each regressor consists of a 2208-channel fully connected (FC) layer with Mish activation [31] and a 256-channel FC layer with sigmoid activation. Batch normalization [20] uses these FC layers. Finally, we predict a camera model by recovering the ranges of the normalized camera parameters to their original ranges. We scale the input images to $224 \times 224$ pixels following conventional studies [17,29,42].

### 3.4. Harmonic non-grid bearing loss

Unlike a loss function based on image reconstruction, Wakai *et al.* proposed the non-grid bearing loss function $L$ [42] based on projecting image coordinates to world coordinates as

$$L_\alpha = \frac{1}{n} \sum_{i=1}^n ||\mathbf{p_i} - \hat{\mathbf{p}}_\mathbf{i}||_2, \tag{4}$$

where $n$ is the number of sampling points; $\alpha$ is a parameter, $\alpha = \{\theta, \psi, f, k_1\}$; and $\hat{\mathbf{p}}$ is the ground-truth value of world coordinates $\mathbf{p}$. Note that $L_\theta$ indicates that the loss function uses a predicted $\theta$ and ground-truth parameters for $\psi$, $f$, and $k_1$. Additionally, $L_\psi$, $L_f$, and $L_{k_1}$ are determined in the same manner. We obtain the world coordinates $\mathbf{p}$ from the image coordinates in sampled points. The sampled points are projected from a unit sphere. For sampling on the unit sphere, we use uniform distribution within valid incident angles that depend on $k_1$. The loss function achieved stably to calibrate cameras using the unit sphere. The joint loss is defined as

$$L = w_\theta L_\theta + w_\psi L_\psi + w_f L_f + w_{k_1} L_{k_1}, \tag{5}$$

where $w_\theta$, $w_\psi$, $w_f$, and $w_{k_1}$ are the joint weights of $\theta$, $\psi$, $f$, and $k_1$, respectively. Although this loss function can effectively train networks, we need to determine the joint weights

for each camera parameter. Wakai *et al.* [42] used the joint weights set to the same values. To determine the optimal joint weights, they needed to repeat training and validation. However, they did not search for the optimal joint weights because of large computational cost.

To address this problem, we surprisingly found that numerical simulations instead of training can analyze the loss landscapes. This loss function can be divided into two steps: predicting camera parameters from a image and projecting sampled points using camera parameters. The latter step requires only the sampled points and camera parameters. Therefore, we focused on the latter step independent of input images. Figure 2 (a) shows the loss landscapes for camera parameters along normalized camera parameters. The landscapes express that the magnitude of loss values of the focal length is the smallest among $\theta$, $\psi$, $f$, and $k_1$, and the focal length is relatively hard to train. Our investigation suggests that the optimal joint loss weights $w$ are estimated as follows: We calculate areas $S$ under the loss function $L$ for $\theta$, $\psi$, $f$, and $k_1$. Assuming practical conditions, we set the ground-truth values to 0.5, which means that the center of the normalized parameter ranges from 0 to 1 in Figure 2 (a). This area $S$ integrating $L$ for the interval 0 to 1 is illustrated in Figure 2 (b) and is given by

$$S_\alpha = \int_0^1 L_\alpha d\alpha = \int_0^1 \frac{1}{n} \sum_{i=1}^n ||\mathbf{p_i} - \hat{\mathbf{p}}_\mathbf{i}||_2 \ d\alpha, \tag{6}$$

These areas $S$ represent the magnitude of each loss for $\theta$, $\psi$, $f$, and $k_1$. Therefore, we define the joint weights $w$ in Equation (5) using normalization as follows:

$$w_\alpha = \tilde{w}_\alpha \ / \ W, \tag{7}$$

where $\tilde{w}_\alpha = 1/S_\alpha$ and $W = \sum_\alpha \tilde{w}_\alpha$. We call a loss function using the weights in Equation (7) "harmonic non-grid bearing loss (HNGBL)." As stated above, our joint weights can alleviate the bias of the magnitude of the loss for camera parameters. Remarkably, we determine these weights without training.

## 4. Experiments

To validate the adaptiveness of our method to various types of fisheye cameras, we conducted massive experiments using large-scale synthetic images and off-the-shelf fisheye cameras.

### 4.1. Datasets

We used two large-scale datasets of outdoor panoramas called the StreetLearn dataset (Manhattan 2019 subset) [30] and the SP360 dataset [6]. First, we divided each dataset into train and test sets following in [42]: $55,599$ train and 161 test images for StreetLearn, and $19,038$ train

**Table 2.** Distribution of the camera parameters for our train set

| Parameters | Distribution | Range or values[1] |
|---|---|---|
| Pan $\phi$ | Uniform | $[0, 360)$ |
| Tilt $\theta$ | Mix<br>Normal<br>Uniform | Normal 70%, Uniform 30%<br>$\mu = 0, \sigma = 15$<br>$[-90, 90]$ |
| Roll $\psi$ | Mix<br>Normal<br>Uniform | Normal 70%, Uniform 30%<br>$\mu = 0, \sigma = 15$<br>$[-90, 90]$ |
| Aspect ratio | Varying | {1/1 9%, 5/4 1%, 4/3 66%,<br>3/2 20%, 16/9 4%} |
| Focal length $f$ | Uniform | $[6, 15]$ |
| Distortion $k_1$ | Uniform | $[-1/6, 1/3]$ |
| Max angle $\eta_{\max}$ | Uniform | $[84, 96]$ |

[1] Units: $\phi, \theta, \psi$, and $\eta_{\max}$ [deg]; $f$ [mm]; $k_1$ [dimensionless]

**Table 3.** Off-the-shelf fisheye cameras with experimental IDs

| ID | Camera body | Camera lens |
|---|---|---|
| 1 | Canon EOS 6D Mark II | Canon EF8-15mm F4L Fisheye USM |
| 2 | Canon EOS 6D Mark II | Canon EF15mm F2.8 Fisheye |
| 3 | Panasonic LUMIX GM1 | Panasonic LUMIX<br>G FISHEYE 8mm F3.5 |
| 4 | FLIR BFLY-U3-23S6C | FIT FI-40 |
| 5 | FLIR FL3-U3-88S2 | FUJIFILM FE185C057HA-1 |
| 6 | KanDao QooCam8K | Built-in |

and 55 test images for SP360. Second, we generated image patches, with a 224-pixel image height ($H_{img}$) and image width ($W_{img} = H_{img} \cdot A$), where $A$ is the aspect ratio, from panorama images: $555,990$ train and $16,100$ test image patches for StreetLearn, and $571,140$ train and $16,500$ test image patches for SP360. Table 2 shows the random distribution of the train set when we generated image patches using camera models with the maximum incident angle $\eta_{\max}$. The test set used the uniform distribution instead of the mixed and varying distribution used for the train set. During the generation step, we set the minimum image circle diameter to the image height, assuming practical conditions.

### 4.2. Off-the-shelf fisheye cameras

We evaluated off-the-shelf fisheye cameras because fisheye cameras have complex lens distortion, unlike narrow-view cameras. Table 3 shows various fisheye cameras that we used for evaluation. Note that we only used the front camera in the QooCam8K camera, which has both front and rear cameras. We captured outdoor fisheye images in Kyoto, Japan using the off-the-shelf cameras.

**Table 4.** Feature summarization of the conventional methods and our method

| Method | DL[2] | Rot[2] | Dist[2] | Projection |
|---|---|---|---|---|
| Santana-Cedrés [37] | | | ✓ | Perspective |
| Liao [24] | ✓ | | ✓ | Perspective |
| Yin [47] | ✓ | | ✓ | Generic camera |
| Chao [7][1] | ✓ | | ✓ | – |
| López-Antequera [29] | ✓ | ✓ | ✓ | Perspective |
| Wakai [42] | ✓ | ✓ | ✓ | Equisolid angle |
| Ours | ✓ | ✓ | ✓ | Proposed generic camera |

[1] Using a generator for undistortion

[2] DL denotes learning-based method; Rot denotes rotation; Dist denotes distortion

### 4.3. Parameter and network settings

To simplify the camera model, we fixed $d_u = d_v$ and the principal points $(c_u, c_v)$ as the center of the image. Because the scale factor depends on the focal length and image sensor size, which is arbitrary for undistortion, we assumed that the image sensor height was 24 mm, which corresponds to a full-size image sensor. We ignored the arbitrary translation vector $\mathbf{t}$. Because the origin of the pan angle is arbitrary, we provided the pan angle for training and evaluation. Therefore, we focused on four trainable parameters, that is, tilt angle $\theta$, roll angle $\psi$, focal length $f$, and a distortion coefficient $k_1$, in our method. Note that we considered camera rotation based on the horizontal line, unlike calibration methods [21, 22] under the Manhattan world assumption.

We optimized our network for a 32 mini-batch size using a rectified Adam optimizer [27], whose weight decay was 0.01. We set the initial learning rate to $1 \times 10^{-4}$ and multiplied the learning rate by 0.1 at the 50th epoch. Additionally, we set the joint weights in Equation (5) using $w_\theta = 0.103$, $w_\psi = 0.135$, $w_f = 0.626$, and $w_{k_1} = 0.136$.

### 4.4. Experimental results

In Table 4, we summarize the features of the conventional methods. We implemented the methods according to the implementation details provided in the corresponding papers, with the exception that StreetLearn and SP360 were used for training the methods of Chao [7], López-Antequera [29], and Wakai [42]. For the method of Santana-Cedrés [37], we excluded test images with few lines because this method requires many lines for calibration.

#### 4.4.1 Parameter and reprojection errors

To validate the accuracy of the predicted camera parameters, we compared methods that can predict rotation and distortion parameters. We evaluated the mean absolute errors of the camera parameters and mean reprojection errors (REPE) on the test set for our generic camera model. Table 5 shows that our method achieved the lowest mean ab-

**Table 5.** Comparison of the absolute parameter errors and reprojection errors on the test set for our generic camera model

| Method | StreetLearn | | | | | SP360 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Mean absolute error ↓ | | | | REPE ↓ | Mean absolute error ↓ | | | | REPE ↓ |
| | Tilt $\theta$ [deg] | Roll $\psi$ [deg] | $f$ [mm] | $k_1$ | [pixel] | Tilt $\theta$ [mm] | Roll $\psi$ [deg] | $f$ [mm] | $k_1$ | [pixel] |
| López-Antequera [29] | 27.60 | 44.90 | 2.32 | – | 81.99 | 28.66 | 44.45 | 3.26 | – | 84.56 |
| Wakai [42] | 10.70 | 14.97 | 2.73 | – | 30.02 | 11.12 | 17.70 | 2.67 | – | 32.01 |
| Ours w/o HNGBL[1] | 7.23 | 7.73 | 0.48 | 0.025 | 12.65 | 6.91 | 8.61 | 0.49 | 0.030 | 12.57 |
| Ours | **4.13** | **5.21** | **0.34** | **0.021** | **7.39** | **3.75** | **5.19** | **0.39** | **0.023** | **7.39** |

[1] "Ours w/o HNGBL" refers to replacing HNGBL with non-grid bearing loss [42]

solute errors and REPE among all methods. This REPE reflected the errors of both extrinsic and intrinsic parameters. To calculate the REPE, we generated $32,400$ uniform world coordinates on the unit sphere within less than $90°$ because of the lack of calibration points for the image-based calibration methods. López-Antequera's method [29] did not seem to work well because it expects non-fisheye input images. Our method substantially reduced focal length errors and camera rotation errors (tilt and roll angles) by $86\%$ and $66\%$, respectively, on average for the two datasets compared with Wakai's method [42]. Furthermore, our method reduces the REPE by $76\%$ on average for the two datasets compared with Wakai's method [42]. Therefore, our method predicted accurate extrinsic and intrinsic camera parameters.

We also evaluated our method, referred to as "Ours w/o HNGBL," replacing our loss function with non-grid bearing loss [42] to analyze the performance of our loss function, as shown in Table 5. This result demonstrates that our loss function effectively reduced the rotation errors in the tilt and roll angles by $3.05°$ on average for the two datasets compared with the "Ours w/o HNGBL" case. In addition to rotation errors, the REPE for our method with HNGBL was 5.22 pixels on average for the two datasets smaller than that for "Ours w/o HNGBL." These results suggest that our loss function enabled networks to accurately predict not only focal length but also other camera parameters.

### 4.4.2 Comparison using PSNR and SSIM

To demonstrate validity and effectiveness in images, we used the peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) [43] for intrinsic parameters. When performing undistortion, extrinsic camera parameters are arbitrary because we consider only intrinsic camera parameters, image coordinates, and incident angles. Table 6 shows the performance of undistortion on the test set for our generic camera model. Our method notably improved the image quality of undistortion by 7.28 for the PSNR and 0.206 for the SSIM on average for the two datasets compared with Wakai's method [42].

To validate the dependency of the four types of fisheye

**Table 6.** Comparison of mean PSNR and SSIM on the test set for our generic camera model

| Method | StreetLearn | | SP360 | |
|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| Santana-Cedrés [37] | 14.65 | 0.341 | 14.26 | 0.390 |
| Liao [24] | 13.71 | 0.362 | 13.85 | 0.404 |
| Yin [47] | 13.91 | 0.349 | 14.03 | 0.390 |
| Chao [7] | 16.13 | 0.409 | 15.88 | 0.449 |
| López-Antequera [29] | 17.88 | 0.499 | 16.24 | 0.486 |
| Wakai [42] | 21.57 | 0.622 | 20.98 | 0.639 |
| Ours w/o HNGBL[1] | 27.41 | 0.801 | 26.49 | 0.801 |
| Ours | **29.01** | **0.838** | **28.10** | **0.835** |

[1] "Ours w/o HNGBL" refers to replacing HNGBL with non-grid bearing loss [42]

camera models, we also evaluated the performance on the trigonometric function models in Table 7. Although orthogonal projection decreased PSNR, our method addressed all the trigonometric function models; hence, our method had the highest PSNR in all cases. This suggests that our generic camera model precisely behaved like a trigonometric function model. Therefore, our method has the potential to calibrate images from various fisheye cameras.

### 4.4.3 Qualitative evaluation

We evaluated the performance of undistortion and full recovery not only for synthetic images but also off-the-shelf cameras to describe the image quality after calibration.
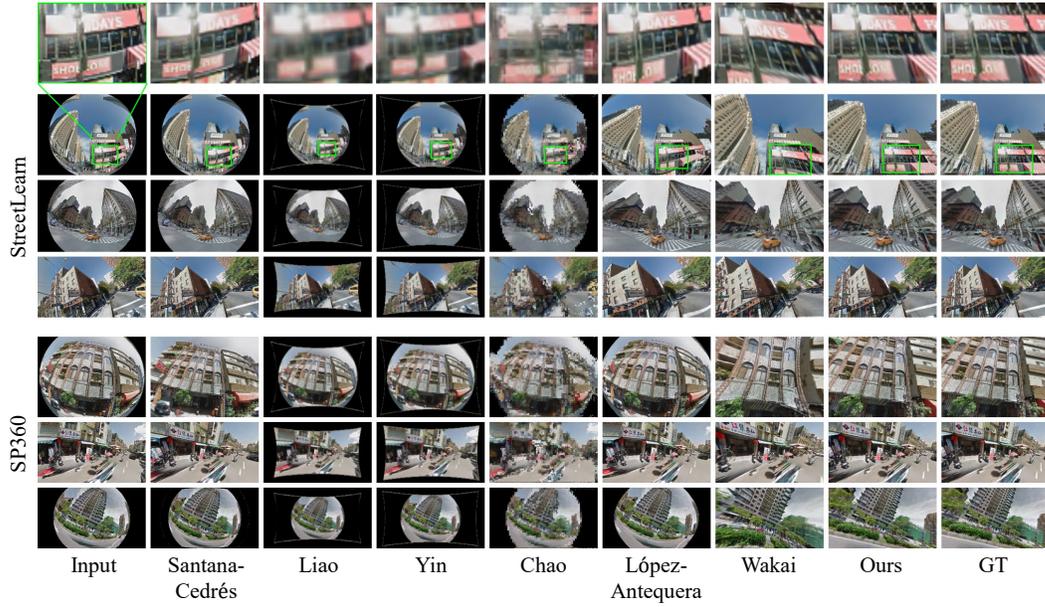
**Synthetic images:** Figure 3 shows the qualitative results on the test set for our generic camera model. Our results are the most similar to the ground-truth images in terms of undistortion, and fully recovering rotation and fisheye distortion. Our method worked well for various types of distortion and scaling. By contrast, it was difficult to calibrate circumferential fisheye images with large distortion using Santana-Cedrés's method [37], Liao's method [24], Yin's method [47], and Chao's method [7]. Furthermore, López-Antequera's [29] and Wakai's [42] methods did not remove distortion, although the scale was close to the ground truth.

When fully recovering rotation and distortion, López-Antequera's [29] and Wakai's [42] methods tended to predict camera rotation with large errors in the tilt and roll an-
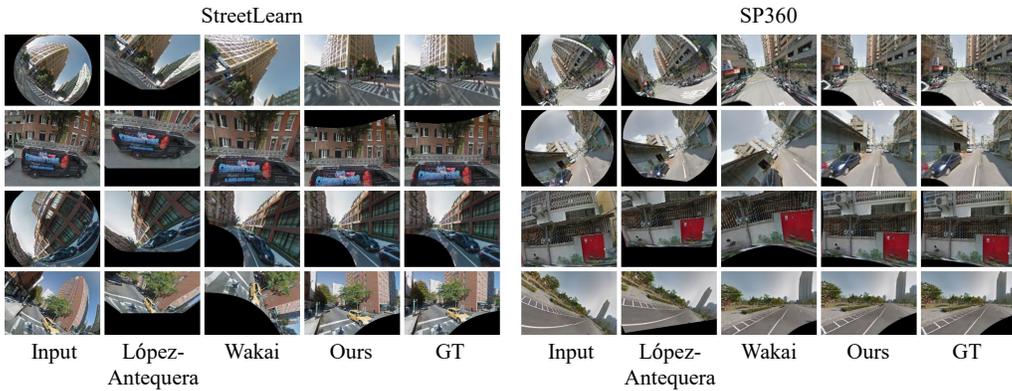
**Table 7.** Comparison of mean PSNR on the test set for the trigonometric function models

| Method | StreetLearn | | | | | SP360 | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | Stereo-graphic | Equi-distance | Equisolid angle | Ortho-gonal | All | Stereo-graphic | Equi-distance | Equisolid angle | Ortho-gonal | All |
| Santana-Cedrés [37] | 14.68 | 13.20 | 12.49 | 10.29 | 12.66 | 14.25 | 12.57 | 11.77 | 9.34 | 11.98 |
| Liao [24] | 13.63 | 13.53 | 13.52 | 13.74 | 13.60 | 13.76 | 13.66 | 13.67 | 13.92 | 13.75 |
| Yin [47] | 13.81 | 13.62 | 13.59 | 13.77 | 13.70 | 13.92 | 13.74 | 13.72 | 13.94 | 13.83 |
| Chao [7] | 15.86 | 15.12 | 14.87 | 14.52 | 15.09 | 15.60 | 15.02 | 14.83 | 14.69 | 15.03 |
| López-Antequera [29] | 17.84 | 16.84 | 16.43 | 15.15 | 16.57 | 15.72 | 14.94 | 14.68 | 14.52 | 14.97 |
| Wakai [42] | 22.39 | 23.62 | 22.91 | 17.79 | 21.68 | 22.29 | 22.65 | 21.79 | 17.54 | 21.07 |
| Ours w/o HNGBL[1] | 26.49 | 29.08 | 28.56 | **23.97** | 27.02 | 25.35 | 28.53 | 28.26 | 23.85 | 26.50 |
| Ours | **26.84** | **30.10** | **29.69** | 23.70 | **27.58** | **25.74** | **29.28** | **28.95** | **23.93** | **26.98** |

[1] "Ours w/o HNGBL" refers to replacing HNGBL with non-grid bearing loss [42]
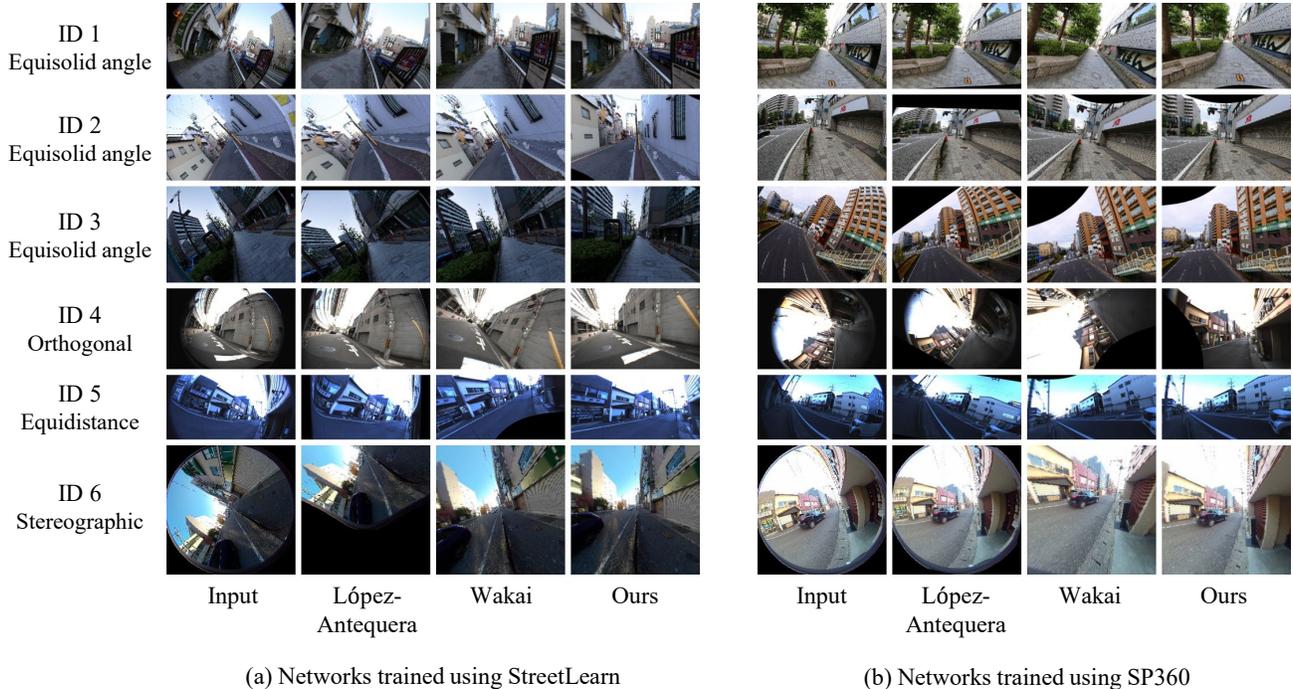


(a) Only undistortion



(b) Fully recovered rotation and distortion

**Figure 3.** Qualitative results on the test images for our generic camera model. (a) Undistortion results shown in the input image, results of the compared methods (Santana-Cedrés [37], Liao [24], Yin [47], Chao [7], López-Antequera [29], and Wakai [42]), our method, and the ground-truth image from left to right. (b) Fully recovered rotation and distortion shown in the input image, results of the compared methods (López-Antequera [29] and Wakai [42]), our method, and the ground-truth image from left to right.

| ID 1<br>Equisolid angle | | | | |
| ID 2<br>Equisolid angle | | | | |
| ID 3<br>Equisolid angle | | | | |
| ID 4<br>Orthogonal | | | | |
| ID 5<br>Equidistance | | | | |
| ID 6<br>Stereographic | | | | |
| | Input | López-<br>Antequera | Wakai | Ours |

(a) Networks trained using StreetLearn

(b) Networks trained using SP360

**Figure 4.** Qualitative results of fully recovering rotation and fisheye distortion for the off-the-shelf cameras shown in the input image, results of the compared methods (López-Antequera [29] and Wakai [42]), and our method from left to right for each image. The IDs correspond to IDs in Table 3, and the projection names are attached to the IDs from specifications (ID: 3–5) and our estimation (ID: 1, 2, and 6). Qualitative results of the methods trained using StreetLearn [30] and SP360 [6] datasets as shown in (a) and (b), respectively.

gles. As shown in Figure 3, our synthetic images consisted of zoom-in images of parts of buildings and zoom-out images of skyscrapers. Our method processed both types of images, that is, it demonstrated scale robustness.

**Off-the-shelf cameras:** We also validated calibration methods using off-the-shelf fisheye cameras to analyze the performance of actual complex fisheye distortion. Note that studies on the conventional learning-based methods in Table 4 reported evaluation results using only synthetic fisheye images. Figure 4 shows the qualitative results of fully recovering rotation and fisheye distortion for methods that predicted intrinsic and extrinsic camera parameters. These methods were trained using the StreetLearn [30] or SP360 [6] datasets. The results for López-Antequera's method had distortion and/or rotation errors. Our method outperformed Wakai's method [42], which often recovered only distortion for all our cameras. Our fully recovered images demonstrated the effectiveness of our method for off-the-shelf fisheye cameras with various types of projection.

In all the calibration methods, images captured by off-the-shelf cameras seemingly degraded the overall performance in the qualitative results compared with synthetic images. This degradation probably occurred because of the complex distortion of off-the-shelf fisheye cameras and

the dataset domain mismatch between the two panorama datasets and our captured images. Overall, our method outperformed the conventional methods in the qualitative evaluation of off-the-shelf cameras. As described above, our method precisely recovered both rotation and fisheye distortion using our generic camera model.

## 5. Conclusion

We proposed a learning-based calibration method using a new generic camera model to address various types of camera projection. Additionally, we introduced a novel loss function that has optimal joint weights determined without training. These weights can alleviate the bias of the magnitude of each loss for four camera parameters. As a result, we enabled networks to precisely predict both extrinsic and intrinsic camera parameters. Extensive experiments demonstrated that our proposed method substantially outperformed conventional geometric-based and learning-based methods on two large-scale datasets. Moreover, we demonstrated that our method fully recovered rotation and distortion using various off-the-shelf fisheye cameras. To improve the calibration performance in off-the-shelf cameras, in future work, we will study the dataset domain mismatch.

# References

[1] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. In *Proceedings of International Conference on Knowledge Discovery and Data Mining (KDD)*, 2019. 2

[2] M. Alemán-Flores, L. Alvarez, L. Gomez, and D. Santana-Cedrés. Automatic lens distortion correction using one-parameter division models. *Image Processing On Line (IPOL)*, 4:327–343, 2014. 2

[3] H. A. Alhaija, S. K. Mustikovela, L. Mescheder, and A. Geiger. Augmented reality meets computer vision: Efficient data generation for urban driving scenes. *International Journal of Computer Vision (IJCV)*, 126(9):961–972, 2018. 1

[4] L. Alvarez, L. Gómez, and J.R. Sendra. An algebraic approach to lens distortion by line rectification. *Journal of Mathematical Imaging and Vision*, 35:36–50, 2009. 3

[5] Z. Cao, T. Simon, S. Wei, and Y. Sheikh. Realtime multi-person 2D pose estimation using part affinity fields. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1302–1310, 2017. 1

[6] S. Chang, C. Chiu, C. Chang, K. Chen, C. Yao, R. Lee, and H. Chu. Generating 360 outdoor panorama dataset with reliable sun position estimation. In *Proceedings of SIGGRAPH Asia*, 2018. 2, 4, 8

[7] C. Chao, P. Hsu, H. Lee, and Y. Wang. Self-supervised deep learning for fisheye image rectification. In *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 2248–2252, 2020. 2, 5, 6, 7

[8] Y. Chen, C. Schmid, and C. Sminchisescu. Self-supervised learning with geometric constraints in monocular video: Connecting flow, depth, and camera. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 7062–7071, 2019. 2

[9] Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2, 3

[10] B. Davidson, M. S. Alvi, and J. F. Henriques. 360$^o$ Camera alignment via segmentation. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 579–595, 2020. 2

[11] Z. Fu, Q. Liu, Z. Fu, and Y. Wang. Template-free visual tracking with space-time memory networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021. 1

[12] D. B. Gennery. Generalized camera calibration including fish-eye lenses. *International Journal of Computer Vision (IJCV)*, 68:239–266, 2006. 3

[13] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2014. 2

[14] I. J. Goodfellow, O. Vinyals, and A. M. Saxe. Qualitatively characterizing neural network. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2015. 2, 3

[15] A. Gordon, H. Li, R. Jonschkowski, and A. Angelova. Depth from videos in the wild: Unsupervised monocular depth learning from unknown cameras. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 8976–8985, 2019. 2

[16] R. Groenendijk, S. Karaoglu, T. Gevers, and T. Mensink. Multi-loss weighting with coefficient of variations. In *Proceedings of IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1469–1478, 2020. 2

[17] Y. Hold-Geoffroy, K. Sunkavalli, J. Eisenmann, M. Fisher, E. Gambaretto, S. Hadap, and J. Lalonde. A perceptual measure for deep single image camera calibration. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2354–2363, 2018. 1, 4

[18] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269, 2017. 3

[19] Z. Huang, Y. Xu, J. Shi, X. Zhou, H. Bao, and G. Zhang. Prior guided dropout for robust visual localization in dynamic environments. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 2791–2800, 2019. 2

[20] S. Ioffe and C. Szegedy. Batch Normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of International Conference on Machine Learning (ICML)*, volume 37, pages 448–456, 2015. 4

[21] J. Lee, H. Go, H. Lee, S. Cho, M. Sung, and J. Kim. CTRL-C: Camera calibration TRansormer with line-classification. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2021. 5

[22] J. Lee, M. Sung, H. Lee, and J. Kim. Neural geometric parser for single image camera calibration. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2020. 5

[23] H. Li, Z. Xu, G. Taylor, C. Studer, and T. Goldstein. Visualizing the loss landscape of neural nets. In *Proceedings of Advances in Neural Information Processing Systems (NeurIPS)*, 2018. 2, 3

[24] K. Liao, C. Lin, and Y. Zhao. A deep ordinal distortion estimation approach for distortion rectification. *IEEE Transactions on Image Processing (TIP)*, 30:3362–3375, 2021. 1, 2, 5, 6, 7

[25] M. Lin, Q. Chen, and S. Yan. Network in network. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2014. 3

[26] C. Liu, L. Chen, F. Schroff, H. Adam, W. Hua, A. L. Yuille, and L. Fei-Fei. Auto-DeepLab: Hierarchical neural architecture search for semantic image segmentation. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 82–92, 2019. 1

[27] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han. On the variance of the adaptive learning rate and beyond. In *Proceedings of International Conference on Learning Representations (ICLR)*, 2020. 5

[28] A. Locher, M. Perdoch, and L. V. Gool. Progressive prioritized multi-view stereo. In *Proceedings of IEEE Conference*

*on Computer Vision and Pattern Recognition (CVPR)*, pages 3244–3252, 2016. 1

[29] M. López-Antequera, R. Marí, P. Gargallo, Y. Kuang, J. Gonzalez-Jimenez, and G. Haro. Deep single image camera calibration with radial distortion. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11809–11817, 2019. 1, 2, 4, 5, 6, 7, 8

[30] P. Mirowski, A. Banki-Horvath, K. Anderson, D. Teplyashin, K. M. Hermann, M. Malinowski, M. K. Grimes, K. Simonyan, K. Kavukcuoglu, A. Zisserman, and R. Hadsell. The StreetLearn environment and dataset. *arXiv preprint arXiv:1903.01292*, 2019. 2, 4, 8

[31] D. Misra. Mish: A self regularized non-monotonic neural activation function. In *Proceedings of British Machine Vision Conference (BMVC)*, 2020. 4

[32] Y. Nie, X. Han, S. Guo, Y. Zheng, J. Chang, and J. J. Zhang. Total3DUnderstanding: Joint layout, object pose and mesh reconstruction for indoor scenes from a single image. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 52–61, 2020. 2

[33] G. V. Puskorius and L. A. Feldkamp. Camera calibration methodology based on a linear perspective transformation error model. In *Proceedings of IEEE International Conference on Robotics and Automation (ICRA)*, pages 1858–1860 vol.3, 1988. 3

[34] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi. You Only Look Once: Unified, real-time object detection. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 779–788, 2016. 1

[35] L. Ren, Y. Song, J. Lu, and J. Zhou. Spatial geometric reasoning for room layout estimation via deep reinforcement learning. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 550–565, 2020. 2

[36] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet large scale visual recognition challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. 3

[37] D. Santana-Cedrés, L. Gomez, M. Alemán-Flores, A. Salgado, J. Esclarín, L. Mazorra, and L. Alvarez. An iterative optimization algorithm for lens distortion correction using two-parameter models. *Image Processing On Line (IPOL)*, 6:326–364, 2016. 2, 5, 6, 7

[38] M. R. U. Saputra, P. Gusmao, Y. Almalioglu, A. Markham, and N. Trigoni. Distilling knowledge from a deep pose regressor network. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 263–272, 2019. 2

[39] L. Sha, J. Hobbs, P. Felsen, X. Wei, P. Lucey, and S. Ganguly. End-to-End camera calibration for broadcast videos. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13624–13633, 2020. 2

[40] S. Shah and J. K. Aggarwal. A simple calibration procedure for fish-eye (high distortion) lens camera. In *Proceedings of IEEE International Conference on Robotics and Automation (ICRA)*, pages 3422–3427 vol.4, 1994. 3

[41] R. Y. Tsai. A versatile camera calibration technique for high-accuracy 3D machine vision metrology using off-the-shelf TV cameras and lenses. *IEEE Journal of Robotics and Automation*, 3(4):323–344, 1987. 1, 2

[42] N. Wakai and T. Yamashita. Deep single fisheye image camera calibration for over 180-degree projection of field of view. In *Proceedings of International Conference on Computer Vision Workshops (ICCVW)*, pages 1174–1183, 2021. 1, 2, 4, 5, 6, 7, 8

[43] Z. Wang and A. C. Bovik. Image quality assessment: From error visibility to structural similarity. *IEEE Transactions on Image Processing (TIP)*, 13(4):600–612, 2004. 6

[44] W. Xian, Z. Li, N. Snavely, M. Fisher, J. Eisenman, and E. Shechtman. UprightNet: Geometry-aware camera orientation estimation from single images. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 9973–9982, 2019. 2

[45] F. Xue, X. Wang, Z. Yan, Q. Wang, J. Wang, and H. Zha. Local supports global: Deep camera relocalization with sequence enhancement. In *Proceedings of International Conference on Computer Vision (ICCV)*, pages 2841–2850, 2019. 2

[46] Z. Xue, N. Xue, G. Xia, and W. Shen. Learning to calibrate straight lines for fisheye image rectification. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1643–1651, 2019. 2

[47] X. Yin, X. Wang, J. Yu, M. Zhang, P. Fua, and D. Tao. FishEyeRecNet: A multi-context collaborative deep network for fisheye image rectification. In *Proceedings of European Conference on Computer Vision (ECCV)*, pages 475–490, 2018. 1, 2, 5, 6, 7

[48] Z. Zhang. A flexible new technique for camera calibration. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 22(11):1330–1334, 2000. 1, 2

[49] H. Żołądek. The topological proof of Abel–Ruffini theorem. *Topological Methods in Nonlinear Analysis*, 16:253–265, 2000. 3