# Neural Lens Modeling

Wenqi Xian[1,*]    Aljaž Božič[2]    Noah Snavely[3]    Christoph Lassner[4]
Meta Reality Labs Research[1,2,4]
Cornell University[1,3]
wx97@cornell.edu[1], aljaz@meta.com[2], snavely@cs.cornell.edu[3], classner@meta.com[4]

## Abstract

*Recent methods for 3D reconstruction and rendering increasingly benefit from end-to-end optimization of the entire image formation process. However, this approach is currently limited: effects of the optical hardware stack and in particular lenses are hard to model in a unified way. This limits the quality that can be achieved for camera calibration and the fidelity of the results of 3D reconstruction. In this paper, we propose NeuroLens, a neural lens model for distortion and vignetting that can be used for point projection and ray casting and can be optimized through both operations. This means that it can (optionally) be used to perform pre-capture calibration using classical calibration targets, and can later be used to perform calibration or refinement during 3D reconstruction, e.g., while optimizing a radiance field. To evaluate the performance of our proposed model, we create a comprehensive dataset assembled from the Lensfun database with a multitude of lenses. Using this and other real-world datasets, we show that the quality of our proposed lens model outperforms standard packages as well as recent approaches while being much easier to use and extend. The model generalizes across many lens types and is trivial to integrate into existing 3D reconstruction and rendering systems. Visit our project website at:* [https://neural-lens.github.io](https://neural-lens.github.io).

## 1. Introduction

Camera calibration is essential for many computer vision applications: it is the crucial component mapping measurements and predictions between images and the real world. This makes calibration a fundamental building block of 3D reconstruction and mapping applications, and of any system that relies on spatial computing, such as autonomous driving or augmented and virtual reality. Whereas camera extrinsics and the parameters of a pinhole model can be easily described and optimized, this often does not hold for other parameters of an optical system and, in particular, lenses. Yet

---

[*]Work done during an internship at RLR.
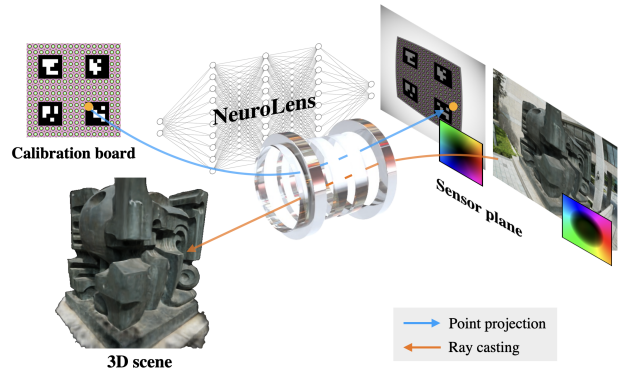[1]The approach is visualized on FisheyeNeRF recordings [23].



Figure 1. **Method Overview.** The optical stack leads to light ray distortion and vignetting. We show that invertible residual networks are a powerful tool to model the distortion for projection and ray casting across many lenses and in many scenarios. Additionally, we propose a novel type of calibration board (top left) that can optionally be used to improve calibration accuracy. For evaluation, we propose the 'SynLens' dataset to evaluate lens models at scale.[1]

lenses have a fundamental influence on the captured image through distortion and vignetting effects.

Recent results in 3D reconstruction and rendering suggest that end-to-end modeling and optimization of the image formation process leads to the highest fidelity scene reproductions [33, 34]. Furthermore, per-pixel gradients are readily available in this scenario and *could* serve as a means to optimize a model of all components of the optical stack to improve reconstruction quality. However, modeling and optimizing lens parameters in full generality and also *differentiably* is hard: camera lenses come in all kinds of forms and shapes (e.g., pinhole, fisheye, catadioptric) with quite different optical effects.

So how can we create a flexible and general and differentiable lens model with enough parameters to approximate any plausible distortion? In classical parametric models, the internals of the camera are assumed to follow a model with a limited number of parameters (usually a polynomial approximation). These approaches work well when the distortion is close to the approximated function, but cannot general-

ize beyond that specific function class. On the other hand, non-parametric models that associate each pixel with a 3D ray have also been explored. These models are designed to model any type of lens, but tend to require dense keypoint measurements due to over-parameterization. Hence, we aim to find models with some level of regularization to prevent such issues, without unnecessarily constraining the complexity of the distortion function. Our key insight is to use an invertible neural network (INN) to model ray distortion, combined with standard camera intrinsics and extrinsics. This means that we model the camera lens as a mapping of two vector fields using a diffeomorphism (*i.e.*, a bijective mapping where both the mapping and its inverse are differentiable), represented by an INN. This approach usefully leverages the invertibility constraints provided by INNs to model the underlying physics of the camera lens.

Our lens model has several advantages. Its formulation makes it easy to differentiate point projection and ray casting operations in deep learning frameworks and it can be integrated into *any* end-to-end differentiable pipeline, with an inductive bias that serves as a useful regularizer for lens models. It is flexible: we can scale the model parameters to adapt to different kinds of lenses.using gradient-based methods for point projection as well as ray casting. This makes our model applicable to pattern-based camera calibration as well as to dense reconstruction where camera parameter refinement is desired. In the case of (optional) marker-based calibration, we suggest to use an end-to-end optimized marker board and keypoint detector. The proposed marker board outperforms several other alternatives in our experiments, and can easily be adjusted to be particularly robust to distortions of different sensor and lens types.

It is currently impossible to evaluate lens models at scale in a standardized way: large-scale camera lens benchmarks including ground truth data simply do not exist. We propose to address this issue by generating a synthetic dataset, called SynLens, consisting of more than 400 different lens profiles from the open-source Lensfun database. To create SynLens, we simulate distortion and vignetting and (optionally) keypoint extraction noise using real lens characteristics to account for a wide variety of lenses and cameras.

We provide qualitative and quantitative comparisons with prior works and show that our method produces more accurate results in a wide range of settings, including precalibration using marker boards, fine-tuning camera models during 3D reconstruction, and using quantitative evaluation on the proposed SynLens dataset. We show that our model achieves subpixel accuracy even with just a few keypoints and is robust to noisy keypoint detections. The proposed method is conceptually simple and flexible, yet achieves state-of-the-art results on calibration problems. We attribute this success to the insight that an INN provides a useful inductive bias for lens modeling and validate this design choise

via ablations on ResNet-based models. To summarize, we claim the following contributions:

- A novel formulation and analysis of an invertible ResNet-based lens distortion model that generalizes across many lens types, is easy to implement and extend;
- A new way to jointly optimize marker and keypoint detectors to increase the robustness of pattern-based calibration;
- A large-scale camera lens benchmark for evaluating the performance of marker detection and camera calibration;
- Integration of the proposed method into a neural rendering pipeline as an example of purely photometric calibration.

## 2. Related Work

**Existing camera calibration methods.** Many 3D computer vision methods assume that lens distortion is radially symmetric around the center of the image. Various camera models such as the radial [13] (bicubic [25]), division [15], FOV models [11], and rational model [8] are used to simulate such radially symmetric distortion. Numerous calibration toolboxes and pipelines [51,52,59] have been developed and integrated to OpenCV [4]. Recently, BabelCalib [32] proposed a robust optimization strategy for parametric models. However, parametric models are only approximate models of real lenses; in practice, the real distortion includes effects caused by complex lens systems (which lead to combinations of different types of distortions) determined by the camera geometry and by the (not perfectly planar) shape of the lens [53].

When calibrating a camera system with an unknown lens it is difficult to decide in advance which particular model fits the real type of camera projection best. To avoid having to choose, one can instead use a single generic model to approximate most common types of projection. A generic camera model [6,18,19,37,43] associates each pixel with a 3D ray. These methods are designed for generality and flexibility and introduce an extreme number of parameters. In practice, classical sparse calibration patterns do not provide enough measurements for such generic models. [2,14] uses these models to obtain dense matches using displays that can encode their pixel positions or interpolate between sparse features. However, interpolation leads to inaccurate and sub-optimal performance. Therefore, models with lower calibration data requirements have been proposed [42]. Recently, Schöps *et al*. [49] extends [42] with a new calibration patterns and detectors to improve the calibration accuracy for generic cameras. [39] replaces the explicit parametric model with a regularization term that forces the underlying distortion map to be smooth.

**Neural network–based camera calibration.** Several prior works treat the optical components of displays and cam-

eras as differentiable layers (neural network layers) that can be trained jointly with the computational blocks of an imaging/display system [20, 50, 54]. Other works estimate camera parameters from single image observations using CNNs [3, 56]. For multi-view, joint optimization of camera parameters and neural scene representations, representative works include BARF [30], NeRF−− [55], Self-Calibrating Neural Radiance Fields [23] and the point-based neural rendering pipeline of Rückert et al. [47]

**Learned markers and keypoint detectors.** Lens models can either be optimized during 3D reconstruction or in a separate calibration stage that uses keypoint positions corresponding to a known 3D structure. Many calibration packages use a checkerboard pattern [5] due to its simplicity and to be able to utilize line fitting to increase corner detection accuracy. Schöps et al. [49] propose a star-based pattern similar to Siemens stars [45] to increase the amount of gradient information available. They use AprilTags [38] to initialize their point search, while we use ArUco tags [16, 46] in a similar way on our proposed marker board.

However, all these boards are manually designed. In contrast, [29] uses a random pattern optimized to produce strong feature responses for keypoint detectors. This leads to significantly more points (on the order of thousands), albeit with lower detection accuracy. Hu et al. [22] propose to use a deep-learning based detector. Grinchuk et al. [17] propose to use a learning-based approach for creating markers by generating binary codes and rendering them on distorted and transformed image patches. Peace et al. [41, 57] use end-to-end trainable systems for marker detection, but focus on fiducial-like markers with a unique marker identification. These systems usually require larger markers with a unique identifier to enable direct estimation of camera pose relative to a single marker. In contrast, we base our board on a marker detector with very high accuracy keypoint detection, as we only care about point detection accuracy and identify points on the board using a few low-accuracy ArUco tags. This leads to a higher number of extracted keypoints and high center point extraction accuracy.

**Invertible Neural Networks.** Our paper models lens distortion using an invertible mapping enforced through the neural network architecture. Invertible neural networks have been studied extensively in the context of normalizing flows, where network inverses are required for computing log-likelihoods for generative models [1, 7, 12, 27, 28]. Since our application does not require the estimation of the Jacobian for generative tasks, we opt to use an invertible residual network due to its expressive power and convergence speed. Invertible residual networks have been applied to many tasks, such as shape deformation [24, 40, 58], image denoising [31], and tone mapping [36]. In this paper, we explore their applicability to the problem of lens distortion.

# 3. Method

The goal of camera calibration is to recover the optimal parameters that describe the camera model at hand given a set of observations. The camera model describes the mapping between points $\mathbf{X} \in \mathbb{R}^3$ in the 3D world and their 2D locations $\mathbf{x} \in \mathbb{R}^2$ on the camera sensor. In this paper, we assume the *projection* component of this mapping to be described by the pinhole camera model. Under this model, the 2D pixel coordinate $x$ can be obtained by:

$$\mathbf{x} = \mathcal{C}(\mathbf{X}) = \text{norm}(\mathbf{K} \cdot (\mathbf{R} \cdot \mathbf{X} + \mathbf{t})), \qquad (1)$$

where $\text{norm}(\mathbf{x}) = (\mathbf{x}[0]/\mathbf{x}[2], \mathbf{x}[1]/\mathbf{x}[2])$, $\mathbf{R}$ and $\mathbf{t}$ are the rotation matrix and translation vector in world-to-camera format, and $\mathbf{K}$ is the intrinsics matrix.

This pinhole model, however, captures only some aspects of the true mapping function for real-world cameras: it assumes that light follows a straight line from the world directly to the sensor plane. This is not the case for real cameras: the optical stack consists of (multiple) lenses with often complex optical properties (e.g., fisheye and catadioptric lenses with wide fields-of-view) that cause visible curvature in the projection of straight lines—the *distortion* component of the mapping. As illustrated in Fig. 2, this non-linear distortion can be modeled by a diffeomorphic function $\mathcal{D}$ that maps ideal coordinates $(u_x, u_y)$ to distorted coordinates $(d_x, d_y)$. As illustrated in Fig. 2, let $u = (u_x, u_y)$ be the normalized coordinates obtained after perspective division but before rescaling by camera intrinsics, the observed pixel coordinate can be obtained by:

$$\mathbf{x} = \text{norm}(\mathbf{K} \cdot \text{hom}(\mathcal{D}(u))) \qquad (2)$$

In contrast to (1) which contains a handful of parameters, our camera model $\mathcal{C}$ contains a bijection which is much more complex to model and $\mathcal{D}$ depends on the physical properties of the camera optics. Hence, it is important to strike a balance between models with sufficiently many parameters that are at the same time constrained to meaningful lens mappings. In our work, we propose to model $\mathcal{D}$ using invertible residual networks and show that they are a strikingly simple, well-suited class of functions for modeling distortions. Using such functions retains the ability to propagate gradients in either the projection (forward) or casting (backward) operation, enabling end-to-end optimization of the camera intrinsics $\mathbf{K}$, extrinsics $\mathbf{R}, \mathbf{T}$, and the distortion mapping $\mathcal{D}$.

In what follows, we will first explain how we parameterize the distortion mapping $\mathcal{D}$ using invertible residual networks in Sec. 3.1. Then, we develop our training objectives in Sec. 3.2, in which we consider a common lens effect of vignetting and incorporate the camera response function (CRF) into our model. Finally, we will describe how to obtain keypoints and their corresponding 3D positions for the case of marker-based calibration. In particular, we propose a
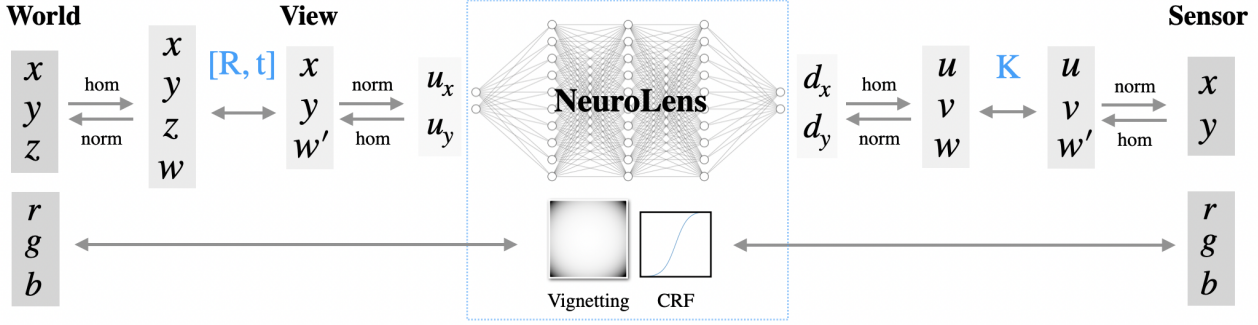
Figure 2. **Mapping Overview.** We illustrate the mapping between three systems: **World** is the world coordinate/color system, **View** is the local camera system, **Sensor** is the sensor coordinate/color system. *hom* and *norm* refer to homogenization and de-homogenization operations, respectively. The figure can be read left-to-right to follow a **projection** operation, right-to-left to follow a **ray casting** operation.

new pattern that enables very accurate keypoint detection in Sec. 3.3.

## 3.1. Camera Distortion Model

As illustrated in Fig. 2, the camera distortion model can be defined as a transformation of a ray from undistorted to distorted directions: $(d_x, d_y) = \mathcal{D}(u_x, u_y)$. In this section, we describe how to parameterize the distortion mapping $\mathcal{D}$.

Intuitively, the distortion transformation can be used in both directions; therefore the process should be *invertible*. Hence we propose to represent the non-linear distortion as a *diffeomorphism*. We can write $\mathcal{D}$ as an invertible function $\mathcal{D} : \mathbb{R}^2 \to \mathbb{R}^2$, where its backward mapping $\mathcal{D}^{-1}$ models the undistortion process. We find that invertible neural networks are a suitable model class for regularizing $\mathcal{D}$ as a smooth, invertible function. Invertible neural networks (INNs) are function approximators that effectively learn differentiable bijections. Networks that are invertible by construction offer a useful advantage: we can train them on a forward mapping and can use the inverse function at no additional cost.

Specifically, we propose to parameterize the distortion mapping $\mathcal{D}$ using Invertible Residual Networks (ResNets), a subclass of INNs introduced by Behrmann *et al.* [1]. Invertible ResNets are composed of residual blocks of the form $f_\theta(x) = x + g_\theta(x)$, where $\theta$ denotes all trainable parameters. Behrmann *et al.* show that $f_\theta$ is invertible if $g_\theta$ is Lipschitz-bounded by 1. In that case, the inverse of $f_\theta$ can be obtained by computing the fixed-point of function $h(x) = y - g_\theta(x)$, where $y$ is the output of $f_\theta(x)$. The fixed-point can be obtained by using the iterative algorithm: $x \leftarrow y - g_\theta(x)$. In practice, we found that a network with width 1024 and four residual blocks is sufficient. For more implementation details, please refer to the supplementary material.

## 3.2. Optimization Objectives

Since we model $\mathcal{D}$ as a differentiable function, it can be used in many optimization scenarios. For instance, it can be used to optimize keypoints and their corresponding 3D positions obtained using calibration targets introduced in

Sec. 3.3, or to optimize all camera parameters together with world model parameters during 3D reconstruction.

**Geometric loss.** A calibration board contains $N$ reference points $\mathbf{X}_i$ whose 3D coordinates are known (in practice, initial estimates for their 3D position can be found using, for example, [60]). The points are assumed to lie in the $XY$-plane, *i.e.*, their $Z$-component is zero. Given a set of 3D-2D points pairs $\{\mathbf{X}_i, \mathbf{x}_i\}_{i=1}^n$, where $\mathbf{x}_i$ is the detected keypoint position in the image, we can minimize the per-view reprojection error:

$$\mathcal{L}_{\text{prj}}(\Theta) = \|\mathcal{C}_\Theta(\mathbf{X}_i) - \mathbf{x}_i\|_2^2, \tag{3}$$

where $\Theta$ includes the camera intrinsics $\mathbf{K}$, extrinsics $\mathbf{R}, \mathbf{T}$, and parameters $\theta$ that define the distortion mapping $\mathcal{D}_\theta$. To improve robustness to outliers, we can optionally apply iterative reweighting to Eq. 3 during optimization.

**Photometric Loss.** Our model can also be optimized using the gradients from a set of 2D image observations, for example as part of a 3D reconstruction. If the camera model describes the image formation procedure correctly, then the color at the pixel location predicted by $\mathcal{C}_\Theta(\mathbf{X})$ should match the color of the actual 3D point $\mathbf{X}$ projected there (assuming constant lighting, exposure and a Lambertian marker material). Suppose $L(\mathbf{X}) \in \mathbb{R}^3$ is the reference color from the calibration board at 3D location $\mathbf{X}$ and $I_i$ is the color of the observed image at pixel location $\mathbf{x}$. Their $\ell_2$ difference can be described by $\|L(\mathbf{X}) - I_i(\mathcal{C}_\Theta(\mathbf{X}))\|_2^2$.

However, this comparison does not take into account optical effects that influence the mapping from radiance to the final image color. Most prominently, we also need to model and estimate vignetting effects (radial falloff) present in many zoom and wide angle lenses [26]. Furthermore, we should take into account the camera response function (CRF), the relationship between the radiance captured by the camera and the resulting sensor readout [9].

To account for such optical effects, we define a function $M(\mathbf{x}, \mathbf{c}) = f(V(\mathbf{x}, \mathbf{c}))$ which takes a pixel location $\mathbf{x}$ and

the incident radiance $\mathbf{c}$ and returns a sensed color taking into account CRF $f$ and vignetting effect $V$. Specifically, we parameterize the vignetting function $V$ by $V_\gamma(\mathbf{x}, \mathbf{c}) = \mathbf{c} \cdot \sigma(\text{interp}(\mathbf{x}, \gamma))$, where $\gamma \in \mathbb{R}^{H \times W}$, interp is bilinear interpolation, and $\sigma$ is a sigmoid function. In our case, we used a fixed CRF $f$ that is known and uniform across the spatial dimensions of the image. We include $\gamma$ as part of the camera parameters $\Theta$, which will be jointly optimized. More general formulations can be used for more complex camera response function and vignetting parameterizations that are appropriate for the camera.

Finally, if keypoints with known radiance are available, for example from the calibration board described in Sec. 3.3, the photometric loss can be used to match the sensed colors to match their expected values. Given $n$ images $\{I_i\}_{i=1}^n$, we can sample color from $m$ points on the calibration board $\{(\mathbf{X}_j, L_j)\}_{j=1}^m$, and define the photometric loss as:

$$\mathcal{L}_{\text{pho}}(\Theta) = \sum_{i,j} \|M(\mathcal{C}_\Theta(\mathbf{X}_j), L_j) - I_i(\mathcal{C}_\Theta(\mathbf{X}_j))\|_2^2. \quad (4)$$

Alternatively, $M$ can be used to map radiance values to color while $\mathcal{C}_\Theta$ are the rays cast for a gradient-based optimization of a radiance field—in that case the gradients can be naturally used to update all relevant parameters (see Sec. 5.4).

### 3.3. Marker-based Calibration

The most common optimization scenario for camera calibration uses an established set of corresponding keypoints to determine $\mathcal{D}$. These are often obtained using a calibration board with a known marker structure that allows for identifying keypoints with high precision. The classical OpenCV [5] library, as well as more recent methods [32,49] use different calibration board types to achieve this. All these board types are hand-designed: their respective patterns yield points with high contrast that can be readily identified. Still, it is not trivial to achieve sub-pixel accurate keypoint detections. In particular, checkerboard corner detection utilizes line-fitting to identify intersection points, and star-shaped pattern detectors rely on symmetric features to identify keypoint centers. All these strategies are non-trivial to implement and are adversely affected by lens distortion.

To address this problem, we propose to optimize the keypoint marker design together with a deep-learning based keypoint detector end-to-end. To represent markers, we create a three-channel tensor that stores an RGB image. To optimize it, we create a simplified model of the image formation process from the marker definition that can contain: printing (small local distortions), lighting (slight intensity changes), motion blur, perspective distortion (viewpoint changes), lens distortion, and color aberration. In our experiments, we use a single marker design optimized for fairly general use by implementing some of the aforementioned effects using
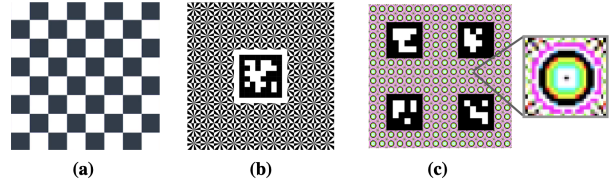


Figure 3. **Keypoint Patterns & Markers. (a)** Checkerboard pattern, **(b)** Star-shaped pattern proposed in [49]. **(c)** Our proposed calibration pattern, allowing for unique localization using ArUco tags [16, 46], and containing high-contrast patterns for accurate keypoint detection. The markers can be optimized specifically for the camera and capture scenario; the size and ratio of markers and tags can be adapted according to the resolution of the camera.

blurring, affine transformations, added noise and color distortion. The detector, a MobileNet-v3 [21] with a simplified 2D location prediction optimized using a Gaussian negative log likelihood of the true keypoint location, has the task of localizing the marker center. This means, we use a fully supervised training for the entire detection process that can be adjusted to match the capture scenario at hand.

The result is empirically superior to other marker shapes and makes better use of color (as shown in Tab. 3): in contrast to the black and white patterns used in manual marker design like checkerboards, our machine-optimized markers use color to maximize cues about keypoint location (see Fig. 3). The symmetry of the marker emerges from the optimization to achieve rotation invariance. The center keypoint is marked black with a small white area around it to maximize contrast and be robust to color bleeding; several circles around it provide additional information to identify and localize it. A pattern board can be readily assembled using these markers by using ArUco [16, 46] markers to identify planar areas and rough sizes, extracting candidate areas and running the pre-trained detector to obtain marker locations. In practice, the confidence prediction from the predicted Gaussian variance helps to filter uncertain detections. Thanks to the high efficiency of MobileNet-v3, the detector runs at multiple frames per second allowing live data acquisition feedback.

## 4. The SynLens Dataset

Evaluating lens models is inherently hard: ground truth is nearly impossible to obtain (since it would require a possibly destructive analysis of equipment), and performing measurements at scale requires a large supply of cameras and lenses. On the other hand, over the last years the LensFun database[2] has steadily grown and accumulated a large set of crowd-sourced high-quality measurements of lens characteristics. Hence, we propose to use it to create a large dataset of high quality *synthetic* lenses that can be used to evaluate calibration models. By creating the data synthetically, we can
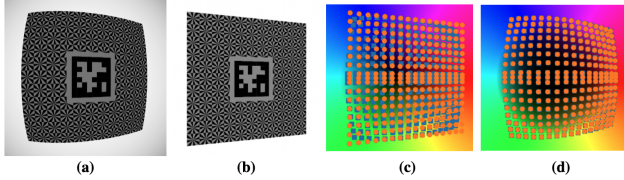
---

[2]https://lensfun.github.io/

Figure 4. **An example from the SynLens Dataset**: **(a)** distorted frame, **(b)** corresponding undistorted and normalized view, **(c)** initialization of keypoints and **(d)** keypoints after optimization. **Blue:** ground truth keypoint positions, **orange:** predicted keypoint positions. **Hue**: offset direction, **saturation**: offset magnitude.

| Models | Formulation ($C$) |
|--------|-------------------|
| Poly3 | $r_d = r_u(1 - k_1 + k_1 r_u^2)$ |
| Poly5 | $r_d = r_u(1 + k_1 r_u^2 + k_2 r_u^4)$ |
| PTLens | $r_d = r_u(a r_u^3 + b r_u^2 + c r_u + 1 - a - b - c)$ |

Table 1. **Analytic equations in LensFun.**

perform calibration in perfect control of noise characteristics and create informative estimates of calibration performance on many consumer devices.

**The Data.** The LensFun database contains more than 3,500 lens models from 40 different camera makers, e.g., Canon, Nikon, action cams, *etc*. For each lens profile, it specifies lens model, focal length, lens distortion, vignetting and chromatic aberration (TCA). High-quality data was collected by photography enthusiasts using the open-source Hugin software. Of this data, we selected 400 lenses by choosing 10 different lens types for each camera maker.

Using this data, we offer dataset users an API to render images, and specifically calibration boards, through these lenses while automatically applying $\mathcal{D}$ and $V$. To test calibration specifically, we provide options to obtain the ground truth positions of projected keypoints. In the following, we describe the API functionality.

**Virtual camera set-up.** We deploy a virtual perspective camera in a synthetic scene using PyTorch3D [44]. It is easy to adjust the virtual camera parameters and to control its pose. We point the camera at a calibration target using several in-plane rotation $\alpha$ and zenith $\beta$ angles. For each scene, we translate the camera off-center and obtain a series of 200 non-parallel images at resolution of $1024 \times 1024$.

**Lens distortion.** From the Lensfun database, lens distortion information is available in one of several predefined formats: PTLens, poly3, poly5 or Adobe Lens (see Tab. 1). According to each calibrated lens model in Lensfun, we synthetically generate distorted and undistorted point pairs in the normalized image domain.

**Vignetting.** The vignette function in the database is parameterized as the polynomial radial loss function $V(r_d) = 1 + k_1 r_d^2 + k_2 r_d^4 + k_3 r_d^6$, where $k1$, $k2$, $k3$ are a set of vi-

gnetting parameters; these model parameters are identical for all color channels. An example of a simulated lens recording a calibration pattern from [49] including distortion and vignetting is shown in Fig. 4. The vignetting effects are clearly visible in (a), whereas (b) shows a successful calibration result. We show recorded and optimized keypoints as well as a visualization of the lens model in subfigures (c) and (d).

## 5. Experiments

In our experiments, we compare the performance of our lens models and marker board on the proposed SynLens dataset with several established methods before presenting results on real-world data for keypoint-based calibration and radiance-field reconstruction on radial and fisheye lenses.

### 5.1. Evaluation on SynLens

| Methods | Camera Models | | | |
|---------|-------|-------|--------|-----|
| | Poly3 | Poly5 | PTLens | Avg |
| Schöps *et al*. [49] | 0.162 | 0.124 | 0.121 | 0.135 |
| Ours | 0.104 | 0.052 | 0.061 | 0.072 |

Table 2. **Reprojection error (RMS) on SynLens by method and lens model for ground truth keypoints.**

| Methods | Keypoint Types | | |
|---------|------------|------|------|
| | Checkboard | Star | Ours |
| OpenCV [5] | 0.152 | 0.175 | 0.129 |
| Schöps *et al*. [49] | 0.178 | 0.141 | 0.158 |
| Ours | 0.154 | 0.130 | 0.114 |

Table 3. **Reprojection error (RMS) on SynLens for detected keypoints by model and keypoint type.**

On SynLens, we establish baseline comparisons with two widely used camera calibration methods and board patterns: (1) the distortion model implemented in OpenCV [5] using all distortion terms, and the board and method from Schöps *et al*. [49], a state-of-the-art generic calibration method with an open-source implementation. Since we evaluate on a synthetic dataset, the ground-truth locations of keypoints can be obtained by transforming them using Eq. 2. We optimize a lens model for each camera lens using the keypoint correspondences; once using the ground-truth project locations of the keypoints, once by using the respective keypoint detector. We then measure the root-mean-squared (RMS) reprojection error (3) on a set of 20 held-out test images uniformly sampled from each sequence.

Tab. 2 shows a breakdown of calibration performance for [49] and our method for the different camera models in the dataset for ground truth keypoint projections. We do not use OpenCV in this table since OpenCV uses the exact same
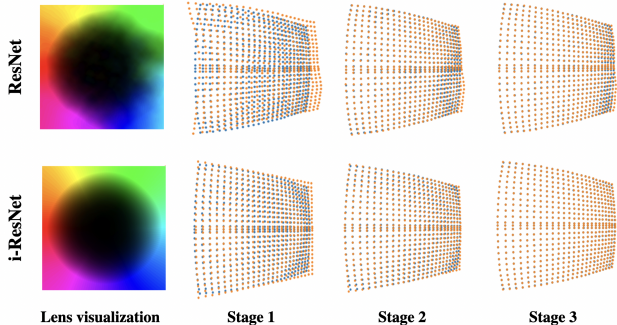
Figure 5. **Comparison between ResNet and invertible ResNet lens models.** Left to right: optimized lens model visualization, and three stages during training, for each: ground-truth keypoints (blue) and projected keypoints (orange) on test images. The rainbow visualization is described in Fig. 4.

parameterization for its model as has been used to generate the data, and therefore unsurprisingly achieves perfect fitting in this scenario. Our method outperforms Schöps *et al.* in this setting for all camera models by a large margin, even though Schöps *et al.* also use a highly parameterized model.

Tab. 3 shows the performance of all methods, using different calibration targets and detection results from the respective keypoint detectors. We expect the best results on the diagonal of the table (each method performing best with its own type of detector pattern and keypoint detector). This mostly holds true, except for OpenCV does better with out marker than with the checkerboard. Our proposed method achieves the overall lowest RMS error with the proposed calibration target in this setting. We analyze how different levels of artificial keypoint noise as well as the severity of the distortion affects the calibration performance of different methods in the supplemental material.

| | i-ResNet | | ResNet | |
|---|---|---|---|---|
| Ratio (tr : val) | Train | Val | Train | Val |
| 1 : 1 | 0.16 | 0.16 | 0.30 | 0.45 |
| 1 : 4 | 0.15 | 0.28 | 0.32 | 0.60 |
| 1 : 8 | 0.15 | 0.33 | 0.37 | 1.56 |

Table 4. **RMS error comparison between ResNet and i-ResNet for different training set sizes.** Total number of keypoints in the validation set remains the same across all experiments.

### 5.2. Comparison with ResNet

The Lipschitz constraint on the invertible ResNet is a powerful regularizer for the proposed model. Compared with standard ResNets, we find that invertible ResNets are less likely to be affected by outliers because they are implicitly constrained to model a smooth function. In Fig. 5, we show a comparison of a ResNet and invertible ResNet trained on a lens with noisy keypoint detections. The ResNet overfits



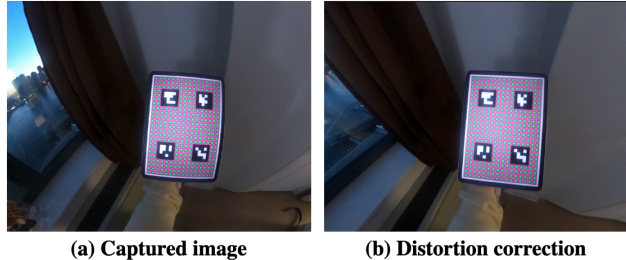**(a) Captured image**   **(b) Distortion correction**

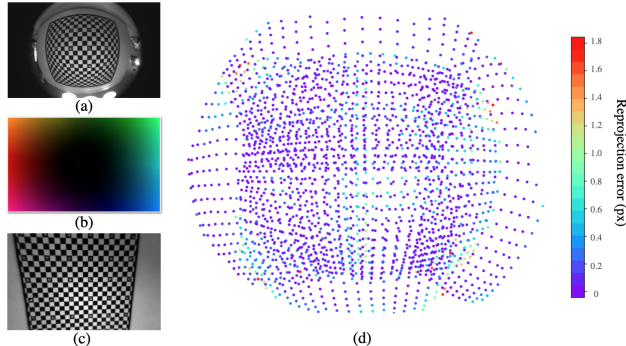Figure 6. **Undistortion of a GoPro super-wide recording.**



Figure 7. **OCamCalib Fisheye camera calibration.** (a) example frame captured by the Fisheye camera, (b) lens distortion map (hue: distortion direction; saturation: distortion magnitude), (c) undistorted image, (d) residuals of reprojected keypoints on test images.

to the noisy measurements present in the training data, for example at the top left corner. In comparison, the invertible ResNet can model accurate lens geometry ofand makes continuous progress towards a reasonable solution over the course of the optimization. In Tab. 4, we show that invertible ResNets are robust to reduced amounts of supervision thanks to their stronger priors.

### 5.3. Evaluation on Real Captures

To ensure that our evaluation results on synthetic data carry over to real-world capture scenarios, we conduct several experiments using challenging wide angle and fisheye lenses. In the first experiment we attempt calibration for a consumer GoPro camera with wide and super-wide lens settings. For data collection, we captured a video of a board with our proposed calibration pattern. We then run keypoint detection and fit our lens model to each camera. Fig. 6 shows the undistortion result for the super-wide lens. For both lenses we achieve slightly better result on super-wide lens setting than OpenCV on held-out test frames: RMS score of OCV 1.50 vs. Ours 1.46, while having comparable results on wide lens settings, OCV 0.56 vs. Ours 0.61.

In the second setting, we extend our experiment to a very challenging scenario: the OCamCalib [48] dataset, with camera field of view ranging from 130° to 266° and the UZH [10] dataset, which consists of eight wide-angle and fisheye cameras with fields of view ranging from 124° to
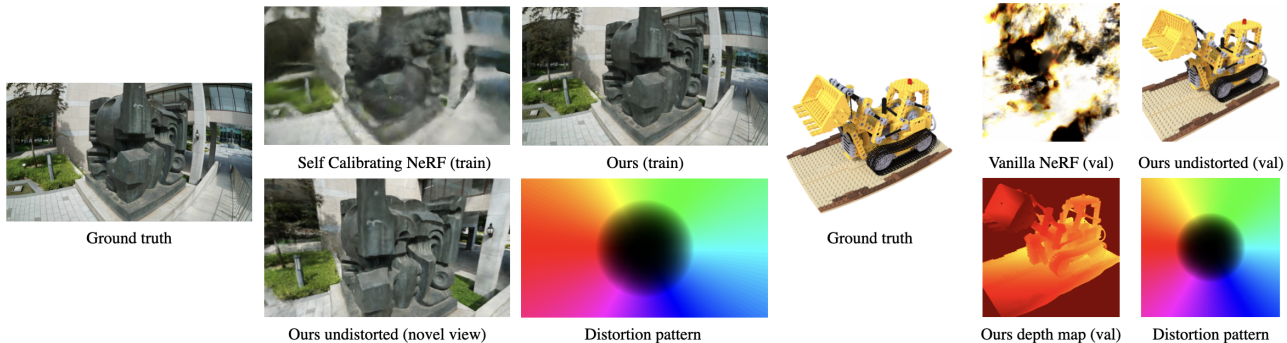
Figure 8. **Neural Radiance Field results.** Qualitative results on the FisheyeNeRF [23] and original NeRF datasets. **Left:** the FisheyeNeRF dataset stretches the capabilities of in-the-wild calibration without keypoint correspondence to the limits. Baseline method SC-NeRF [23], to the best of our knowledge, only shows results on training views, for which our model fares remarkably well. **Right:** Results on a scene from the NeRF dataset. We used a Blender scene and added significant barrel distortion (as visible in the "Ground truth" setting). NeRF [34] fails to reconstruct the scene. Our method manages to retrieve the lens parameters well, resulting in high-quality reconstruction and depth.

$166°$. Keypoint detections are available from a planar chessboard target marked with AprilTags. We compare our results with the state-of-the-art camera calibration framework BabelCalib [32]. As shown in Tab. 5, our method outperforms BabelCalib on most cameras from the UZH dataset. We also visualize the residuals of the reprojected keypoints of test images in Fig. 7 from OCamCalib, on which our method achieves an *unweighted* reprojection error of 0.91 (all points contribute equally to the error metric). This is a comparable score with the BabelCalib system with a significantly simpler model on this challenging data.

| | UZH-DAVIS | | | | UZH-Snapdragon | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Methods | I-1 | I-2 | O-1 | O-2 | I-1 | I-2 | O-1 | O-2 | Mean |
| BabelCalib | 0.31 | 0.69 | 1.58 | 0.44 | **0.28** | 1.10 | 0.68 | **0.28** | 0.67 |
| Ours | **0.25** | **0.52** | **0.49** | **0.36** | 0.64 | **0.57** | **0.34** | 1.08 | **0.53** |

Table 5. **RMS score comparison between BabelCalib and our method on UZH camera dataset.**

### 5.4. Neural Radiance Fields

A significant advantage of our proposed lens model is its two-way differentiability, making it simple to deploy in 3D reconstruction workflows. Neural radiance fields (NeRFs) [34] are a state-of-the-art approach for novel-view synthesis. They optimize a scene model directly from RGB images, given the camera intrinsics and poses. While NeRF achieves high-quality novel views, it requires accurate camera parameters, which can be difficult to obtain in practice, particular for lens parameters, which often require an additional calibration stage. We integrate our neural lens model into a neural rendering framework [35] such that the camera poses, pinhole intrinsics and lens distortion are optimized together with the appearance model, given only RGB observations. The camera intrinsic and extrinsic parameters are initialized using values obtained from a photogrammetry software package, Metashape, yet undistorted. As can be seen in Fig. 8, our approach achieves a high-quality repre-

sentation of camera views and successfully recovers the lens distortion, even in the case of extremely distorted recordings from the FisheyeNeRF dataset [23].

To experiment with significant distortion, but still in a non-fisheye setting, we use the NeRF dataset and augment a Blender scene with barrel distortion. The NeRF reconstruction fails completely in this setting; augmented with our proposed model it succeeds in reconstruction and undistortion without any other changes to the training pipeline.

## 6. Limitations and Future Work

In cases of very extreme lens distortion, it could be helpful to initialize the model with a prior expectation as opposed to starting from an identity initialization. This could help the convergence rate as well as lead to even better solutions. Incorporating lens priors for specific models could also be used for model regularization if that's desired for the specific application, though we found the proposed model to be very stable and usually not needing additional regularizers.

## 7. Conclusion

In this paper, we presented a novel approach for neural lens modelling with a focus on end-to-end optimization, generality and ease-of-use in existing deep learning pipelines. It includes distortion as well as vignetting effects and, thanks to being based on invertible residual network models, can be optimized for projection and raycasting. The model can directly be used to improve the results for 3D reconstructions for radiance field models with hardly any changes to existing implementations. We also introduced an end-to-end differentiable marker-board and point detector that can be used to perform offline calibration. Using our proposed synthetic lens dataset as well as results on GoPro and fisheye camera, we showed that the proposed model generalizes across lenses, cameras and applications and can be a reliable calibration component for future applications and research.

# References

[1] Jens Behrmann, Will Grathwohl, Ricky T. Q. Chen, David Duvenaud, and Joern-Henrik Jacobsen. Invertible residual networks. In Kamalika Chaudhuri and Ruslan Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 573–582. PMLR, 09–15 Jun 2019. 3, 4

[2] Filippo Bergamasco, Luca Cosmo, Andrea Gasparetto, Andrea Albarelli, and Andrea Torsello. Parameter-free lens distortion calibration of central cameras. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3847–3855, 2017. 2

[3] Oleksandr Bogdan, Viktor Eckstein, Francois Rameau, and Jean-Charles Bazin. Deepcalib: A deep learning approach for automatic intrinsic calibration of wide field-of-view cameras. In *Proceedings of the 15th ACM SIGGRAPH European Conference on Visual Media Production*, pages 1–10, 2018. 3

[4] Gary Bradski. The opencv library. *Dr. Dobb's Journal: Software Tools for the Professional Programmer*, 25(11):120–123, 2000. 2

[5] G. Bradski. The OpenCV Library. *Dr. Dobb's Journal of Software Tools*, 2000. 3, 5, 6

[6] Federico Camposeco, Torsten Sattler, and Marc Pollefeys. Non-parametric structure-based calibration of radially symmetric cameras. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2192–2200, 2015. 2

[7] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018. 3

[8] David Claus and Andrew W Fitzgibbon. A rational function lens distortion model for general cameras. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 213–219. IEEE, 2005. 2

[9] Paul E Debevec and Jitendra Malik. Recovering high dynamic range radiance maps from photographs. In *ACM SIGGRAPH 2008 classes*, pages 1–10. ACM, 2008. 4

[10] Jeffrey Delmerico, Titus Cieslewski, Henri Rebecq, Matthias Faessler, and Davide Scaramuzza. Are we ready for autonomous drone racing? the uzh-fpv drone racing dataset. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 6713–6719. IEEE, 2019. 7

[11] Frederic Devernay and Olivier Faugeras. Straight lines have to be straight. *Machine vision and applications*, 13(1):14–24, 2001. 2

[12] Laurent Dinh, David Krueger, and Yoshua Bengio. Nice: Nonlinear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014. 3

[13] C Brown Duane. Close-range camera calibration. *Photogramm. Eng*, 37(8):855–866, 1971. 2

[14] Aubrey K Dunne, John Mallon, and Paul F Whelan. Efficient generic calibration method for general cameras with single centre of projection. *Computer Vision and Image Understanding*, 114(2):220–233, 2010. 2

[15] Andrew W Fitzgibbon. Simultaneous linear estimation of multiple view geometry and lens distortion. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 1, pages I–I. IEEE, 2001. 2

[16] S. Garrido-Jurado, R. Muñoz-Salinas, F. J. Madrid-Cuevas, and R. Medina-Carnicer. Generation of fiducial marker dictionaries using Mixed Integer Linear Programming. *Pattern Recognition*, 51:481–491, mar 2016. 3, 5

[17] Oleg Grinchuk, Vadim Lebedev, and Victor Lempitsky. Learnable visual markers. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016. 3

[18] Michael D Grossberg and Shree K Nayar. A general imaging model and a method for finding its parameters. In *Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001*, volume 2, pages 108–115. IEEE, 2001. 2

[19] Richard Hartley and Sing Bing Kang. Parameter-free radial distortion correction with center of distortion estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 29(8):1309–1321, 2007. 2

[20] Felix Heide, Markus Steinberger, Yun-Ta Tsai, Mushfiqur Rouf, Dawid Pajak, Dikpal Reddy, Orazio Gallo, Jing Liu, Wolfgang Heidrich, Karen Egiazarian, et al. Flexisp: A flexible camera image processing framework. *ACM Transactions on Graphics (ToG)*, 33(6):1–13, 2014. 3

[21] Andrew Howard, Mark Sandler, Bo Chen, Weijun Wang, Liang-Chieh Chen, Mingxing Tan, Grace Chu, Vijay Vasudevan, Yukun Zhu, Ruoming Pang, Hartwig Adam, and Quoc Le. Searching for mobilenetv3. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1314–1324, 2019. 5

[22] Danying Hu, Daniel DeTone, and Tomasz Malisiewicz. Deep charuco: Dark charuco marker pose estimation. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8428–8436, 2019. 3

[23] Yoonwoo Jeong, Seokjun Ahn, Christopher Choy, Anima Anandkumar, Minsu Cho, and Jaesik Park. Self-calibrating neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5846–5854, 2021. 1, 3, 8

[24] Chiyu Jiang, Jingwei Huang, Andrea Tagliasacchi, and Leonidas J Guibas. Shapeflow: Learnable deformation flows among 3d shapes. *Advances in Neural Information Processing Systems*, 33:9745–9757, 2020. 3

[25] Einari Kilpelä. Compensation of systematic errors of image and model coordinates. *Photogrammetria*, 37(1):15–44, 1981. 2

[26] Seon Joo Kim and Marc Pollefeys. Robust radiometric calibration and vignetting correction. *IEEE transactions on pattern analysis and machine intelligence*, 30(4):562–576, 2008. 4

[27] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018. 3

[28] Ivan Kobyzev, Simon JD Prince, and Marcus A Brubaker. Normalizing flows: An introduction and review of current

methods. *IEEE transactions on pattern analysis and machine intelligence*, 43(11):3964–3979, 2020. 3

[29] Bo Li, Lionel Heng, Kevin Koser, and Marc Pollefeys. A multiple-camera system calibration toolbox using a feature descriptor-based calibration pattern. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 1301–1307, 2013. 3

[30] Chen-Hsuan Lin, Wei-Chiu Ma, Antonio Torralba, and Simon Lucey. Barf: Bundle-adjusting neural radiance fields. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5741–5751, 2021. 3

[31] Yang Liu, Zhenyue Qin, Saeed Anwar, Pan Ji, Dongwoo Kim, Sabrina Caldwell, and Tom Gedeon. Invertible denoising network: A light solution for real noise removal. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13365–13374, 2021. 3

[32] Yaroslava Lochman, Kostiantyn Liepieshov, Jianhui Chen, Michal Perdoch, Christopher Zach, and James Pritts. Babelcalib: A universal approach to calibrating central cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15253–15262, 2021. 2, 5, 8

[33] Stephen Lombardi, Tomas Simon, Jason Saragih, Gabriel Schwartz, Andreas Lehrmann, and Yaser Sheikh. Neural volumes: Learning dynamic renderable volumes from images. *ACM Trans. Graph.*, 38(4):65:1–65:14, July 2019. 1

[34] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 1, 8

[35] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *arXiv preprint arXiv:2201.05989*, 2022. 8

[36] Aamir Mustafa, Param Hanji, and Rafal Mantiuk. Distilling style from image pairs for global forward and inverse tone mapping. In *Proceedings of the 19th ACM SIGGRAPH European Conference on Visual Media Production*, pages 1–10, 2022. 3

[37] David Nistér, Henrik Stewénius, and Etienne Grossmann. Non-parametric self-calibration. In *Tenth IEEE International Conference on Computer Vision (ICCV'05) Volume 1*, volume 1, pages 120–127. IEEE, 2005. 2

[38] Edwin Olson. AprilTag: A robust and flexible visual fiducial system. In *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pages 3400–3407. IEEE, May 2011. 3

[39] Linfei Pan, Marc Pollefeys, and Viktor Larsson. Camera pose estimation using implicit distortion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12819–12828, 2022. 2

[40] Despoina Paschalidou, Angelos Katharopoulos, Andreas Geiger, and Sanja Fidler. Neural parts: Learning expressive 3d shape abstractions with invertible neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3204–3215, 2021. 3

[41] J. Brennan Peace, Eric Psota, Yanfeng Liu, and Lance C. Pérez. E2etag: An end-to-end trainable method for generating

and detecting fiducial markers. *The 31st British Machine Vision Conference (BMVC)*, 2020. 3

[42] Srikumar Ramalingam and Peter Sturm. A unifying model for camera calibration. *IEEE transactions on pattern analysis and machine intelligence*, 39(7):1309–1319, 2016. 2

[43] Srikumar Ramalingam, Peter Sturm, and Suresh K Lodha. Towards complete generic camera calibration. In *CVPR*, volume 1, pages 1093–1098. IEEE, 2005. 2

[44] Nikhila Ravi, Jeremy Reizenstein, David Novotny, Taylor Gordon, Wan-Yen Lo, Justin Johnson, and Georgia Gkioxari. Accelerating 3d deep learning with pytorch3d. *arXiv:2007.08501*, 2020. 6

[45] R. Reulke, S. Becker, N. Haala, and U. Tempelmann. Determination and improvement of spatial resolution of the ccd-line-scanner system ads40. *ISPRS Journal of Photogrammetry and Remote Sensing*, 60(2):81–90, 2006. 3

[46] Francisco J. Romero-Ramirez, Rafael Muñoz-Salinas, and Rafael Medina-Carnicer. Speeded up detection of squared fiducial markers. *Image and Vision Computing*, 76:38–47, aug 2018. 3, 5

[47] Darius Rückert, Linus Franke, and Marc Stamminger. Adop: Approximate differentiable one-pixel point rendering. *ACM Transactions on Graphics (TOG)*, 41(4):1–14, 2022. 3

[48] Davide Scaramuzza, Agostino Martinelli, and Roland Siegwart. A flexible technique for accurate omnidirectional camera calibration and structure from motion. In *Fourth IEEE International Conference on Computer Vision Systems (ICVS'06)*, pages 45–45. IEEE, 2006. 7

[49] Thomas Schops, Viktor Larsson, Marc Pollefeys, and Torsten Sattler. Why having 10,000 parameters in your camera model is better than twelve. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2535–2544, 2020. 2, 3, 5, 6

[50] Vincent Sitzmann, Steven Diamond, Yifan Peng, Xiong Dun, Stephen Boyd, Wolfgang Heidrich, Felix Heide, and Gordon Wetzstein. End-to-end optimization of optics and image processing for achromatic extended depth of field and super-resolution imaging. *ACM Transactions on Graphics (TOG)*, 37(4):1–13, 2018. 3

[51] Peter Sturm and Srikumar Ramalingam. A generic concept for camera calibration. In *European Conference on Computer Vision*, pages 1–13. Springer, 2004. 2

[52] Rahul Swaninathan, Michael D Grossberg, and Shree K Nayar. A perspective on distortions. In *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, volume 2, pages II–594. IEEE, 2003. 2

[53] Zhongwei Tang, Rafael Grompone Von Gioi, Pascal Monasse, and Jean-Michel Morel. A precision analysis of camera distortion models. *IEEE Transactions on Image Processing*, 26(6):2694–2704, 2017. 2

[54] Ethan Tseng, Felix Yu, Yuting Yang, Fahim Mannan, Karl ST Arnaud, Derek Nowrouzezahrai, Jean-François Lalonde, and Felix Heide. Hyperparameter optimization in black-box image processing using differentiable proxies. *ACM Trans. Graph.*, 38(4):27–1, 2019. 3

[55] Zirui Wang, Shangzhe Wu, Weidi Xie, Min Chen, and Victor Adrian Prisacariu. Nerf–: Neural radiance fields without

known camera parameters. *arXiv preprint arXiv:2102.07064*, 2021. 3

[56] Wenqi Xian, Zhengqi Li, Matthew Fisher, Jonathan Eisenmann, Eli Shechtman, and Noah Snavely. Uprightnet: geometry-aware camera orientation estimation from single images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9974–9983, 2019. 3

[57] Mustafa B. Yaldiz, Andreas Meuleman, Hyeonjoong Jang, Hyunho Ha, and Min H. Kim. Deepformabletag: End-to-end generation and recognition of deformable fiducial markers. *ACM Trans. Graph.*, 40(4), jul 2021. 3

[58] Guandao Yang, Serge Belongie, Bharath Hariharan, and Vladlen Koltun. Geometry processing with neural fields. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 3

[59] Zhengyou Zhang. A flexible new technique for camera calibration. *IEEE Transactions on pattern analysis and machine intelligence*, 22(11):1330–1334, 2000. 2

[60] Z. Zhang. A flexible new technique for camera calibration. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(11):1330–1334, 2000. 4