# Robust Gyroscope-Aided Camera Self-Calibration

Santiago Cortés Reina
Department of Computer Science
Aalto University
Helsinki, Finland
santiago.cortesreina@aalto.fi

Arno Solin
Department of Computer Science
Aalto University
Helsinki, Finland
arno.solin@aalto.fi

Juho Kannala
Department of Computer Science
Aalto University
Helsinki, Finland
juho.kannala@aalto.fi

*Abstract*—Camera calibration for estimating the intrinsic parameters and lens distortion is a prerequisite for various monocular vision applications including feature tracking and video stabilization. This application paper proposes a model for estimating the parameters on the fly by fusing gyroscope and camera data, both readily available in modern day smartphones. The model is based on joint estimation of visual feature positions, camera parameters, and the camera pose, the movement of which is assumed to follow the movement predicted by the gyroscope. Our model assumes the camera movement to be free, but continuous and differentiable, and individual features are assumed to stay stationary. The estimation is performed online using an extended Kalman filter, and it is shown to outperform existing methods in robustness and insensitivity to initialization. We demonstrate the method using simulated data and empirical data from an iPad.

## I. INTRODUCTION

The growth in the market of smartphones and tablets has brought monocular cameras to even the cheapest of smart-devices. Simultaneously, the improved computational capabilities have made it possible to expand their use into new fields and use cases, such as video calls, payment verification, and augmented reality. However, employing the device camera in pose estimation, video stabilization, or feature tracking requires the camera calibration parameters to be known or estimated.

On smart-devices, the use cases often provide an explicit calibration step that includes capturing a pre-defined calibration pattern from different positions. This procedure is the base for traditional camera calibration (see, *e.g.*, [9] for an overview). Self-calibration is the problem of modeling the internal parameters of a camera (projection matrix and distortion coefficients) without using any known pattern. Luckily there are usually additional sensors available on the devices, typically a low-cost MEMS inertial measurement unit (IMU). The IMU typically provides fast-sampled (up to some hundreds of Hz) readings of the specific force (accelerometer) and turn-rate (gyroscope) in the device's coordinate frame (see, *e.g.*, [26] for a more thorough introduction).

In this work we propose a method for camera self-calibration using the information from a gyroscope rigidly attached to the camera. We use a structure from motion (SfM) approach, where the camera model minimizes the reprojection error over a set of images. The proposed method jointly estimates the camera pose, tracked feature positions, and camera
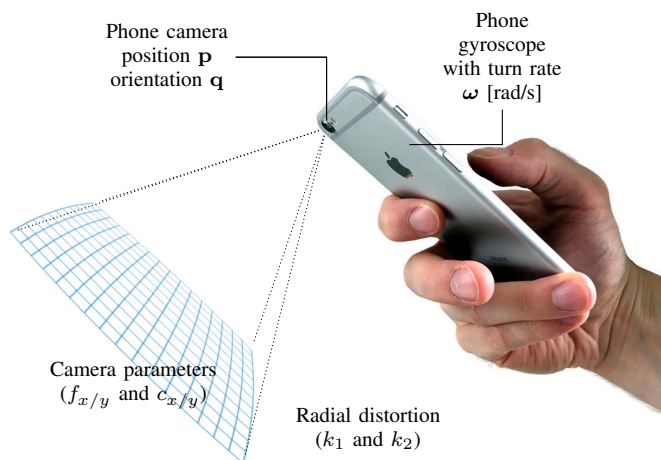


Fig. 1. Illustration of the model setup: A pinhole camera with position $\mathbf{p}$ and orientation quaternion $\mathbf{q}$, auxiliary turn rate $\boldsymbol{\omega}$ from the gyroscope. A set of (unknown) feature locations $\{\mathbf{z}_i\}$ are observed through a camera model with linear and non-linear distortion parameters.

parameters. The method works online employing a state-space formulation, where the fast-sampled gyroscope data is used for forward-predicting the relative camera and feature movement, and the visual data tracking results are then matched to the predictions in a probabilistic fashion. The inference itself is solved by an extended Kalman filter.

Due to its practical importance, camera self-calibration—or auto-calibration—has been studied extensively during the past decades. Outside smartphones, practical applications are found, for example, in traffic surveillance [1], projector camera calibration [18, 29] and robotics for autonomous vehicles [5, 7, 8, 10, 21, 28]. The seminal paper by Faugeras *et al.* [6] introduced the concept of calibrating the intrinsic and extrinsic parameters without a known calibration object or pattern. Since then many methods have been introduced, building upon special kind of motion (*e.g.* purely rotating or planar), special scene geometry (*e.g.* planar scenes or special depth structure, see [11] for discussion), or auxiliary sensors.

For example, Civera *et al.* [4] proposed a sum-of-Gaussians filter approach building upon a SfM approach where there has to be sufficient translation. Additional sensors ease the motion constraints. Using both a gyroscope and an accelerometer for calibration, information can be acquired about the

relative rotation, absolute scale, and the world coordinate frame orientation (see, *e.g.*, [19]). Similar setups are often employed in visual-inertial odometry methods (see, *e.g.*, the discussion in [24, 27]), where the same set of sensors are used and the camera calibration is estimated as a part of the inference. Gyroscopes are also used in video stabilization, and [22] proposed a model for scaling, time offset, and relative pose calibration using the gyroscope. They, however, do not estimate the camera calibration parameters.

There are also other methods which only use the camera and a gyroscope. Karpenko *et al.* [16] proposed a method for calibration which uses a gyroscope together with a quick shake when pointing to far-away objects, while Hwangbo *et al.* [12] constrain the estimation by assuming pure rotation.

Within methods designed for fusing camera and gyroscope data, we are only aware of one previously published approach, which works in the same setting as the proposed method. The method by Jia and Evans was first published in [14] and later refined in [15]. Their method together with ours does not assume special movement and only require visual and gyroscope data.

The proposed model is based on the following assumptions: (i) The camera movement is free, but continuous and differentiable (rather smooth), (ii) The camera rotations follow the gyroscope turn rate observations, (iii) The world coordinates of individual features are unknown, but assumed constant (individual features assumed to stay stationary). In comaprison to [15], this method aims at providing higher robustness to initialization and feature-poor enivironments with only a few visual features being tracked. Where [15] uses a lot of short (two-frame) connections of features, our method uses a few long chains of connected features.

This paper is structured as follows. Section II presents the theoretical methodology in detail starting from the camera calibration model, the proposed state-space model, and finally how to jointly infer all the unknown parameters. The Results section shows the method employed both in a simulation study and a on real-world data. Finally, the assumptions and modeling problems and some future work are discussed.

## II. METHODS

We start by defining the notation used in the camera calibration model, which will later be used in specification of the measurement model. Then we proceed to setting up the state-space model for the gyroscope-aided dynamics of the camera pose, feature positions, and camera parameters. Finally, we couple the forward dynamics with the image data in a visual measurement model. A sketch of the information flow in the method is shown in Figure 2.

### A. Camera Model

We build upon the well-established theory of monocular camera calibration models. A detailed and extensive description of camera models can be found in [9]. Through the camera model, points $(x, y, z)$ in three-dimensional space are projected to image plane coordinates $(u, v)$ using a pinhole
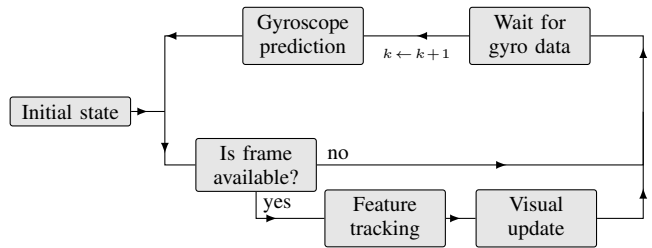


Fig. 2. The information flow in the self-calibration method. The gyroscope prediction loop runs at 100 Hz, and visual updates occur once a new frame is acquired (typically at 10 Hz).

camera model. First the points are rigidly transformed into the camera reference frame, that is origin at the optical center, the $z$-axis perpendicular to the image plane and $y$-axis parallel to the vertical axis of the image. Then the points are projected into the image plane using the pinhole projection. In matrix form

$$\begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \mathbf{K}\,\mathbf{E} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix},$$ (1)

where the intrinsic matrix $\mathbf{K} \in \mathbb{R}^{3 \times 3}$ and extrinsic matrix $\mathbf{E} \in \mathbb{R}^{3 \times 4}$ are parametrized as follows:

$$\mathbf{K} = \begin{pmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \quad \text{and} \quad \mathbf{E} = \begin{pmatrix} \mathbf{R}^\mathsf{T} & -\mathbf{R}^\mathsf{T}\mathbf{p} \end{pmatrix}.$$ (2)

The parameters are the focal length ($f_x$ and $f_y$, separate between dimensions to account for non-square pixels), the origin of the image plane $(c_x, c_y)$, and the position $\mathbf{p} \in \mathbb{R}^3$ and orientation $\mathbf{R} \in \mathbb{R}^{3 \times 3}$ of the camera in the global frame of reference.

The projection operation is entirely linear (in homogeneous space), but real-world lenses usually introduce non-linear mappings. These are known as distortions, the most common distortion is rotationally symmetric around the image center and is modeled by so-called radial distortion coefficients $k_1$ and $k_2$ as follows:

$$\begin{pmatrix} u' \\ v' \end{pmatrix} = \begin{pmatrix} u \\ v \end{pmatrix} (1 + k_1\,r^2 + k_2\,r^4),$$ (3)

where the radial component is determined by

$$r = \sqrt{\left(\frac{u - c_x}{f_x}\right)^2 + \left(\frac{v - c_y}{f_y}\right)^2}.$$ (4)

This camera model is shown in the illustrative sketch in Figure 1, where also the resulting distortion effect is visible. We refer to the non-linear distortion function as 'distort' later in the paper.

### B. State Estimation

In order to include time-dependency between observed frames and rotational information from the gyroscope, we define a state-space model. The model describes a device with a

monocular camera and a rigidly attached gyroscope with a known relative orientation between them (see Fig. 1). An example of such a device is a modern smartphone. The following section shows the non-linear filer designed for information fusion, the non-linear filter allows for estimation and fusion while keeping track of the uncertainty and dependencies between the non-deterministic processes describing the evolution of the calibration parameters, camera pose, and feature locations.

We define a state-space model (see, *e.g.*, [23]), where the state vector is

$$\mathbf{x} = \begin{pmatrix} \mathbf{c}^\mathsf{T} & \mathbf{p}^\mathsf{T} & \mathbf{v}^\mathsf{T} & \mathbf{q}^\mathsf{T} & \mathbf{z}^\mathsf{T} \end{pmatrix}^\mathsf{T}. \tag{5}$$

The variable $\mathbf{c} = (f_x, f_y, c_x, c_y, k_1, k_2)$ contains the internal camera parameters, $\mathbf{p} \in \mathbb{R}^3$ and $\mathbf{v} \in \mathbb{R}^3$ contain the position and velocity of the camera, $\mathbf{q}$ contains the quaternion encoding the orientation of the camera, and $\mathbf{z} \in \mathbb{R}^{3d}$ contains the locations of the features ($d$ is the number of features being tracked).

The non-linear state-space model with the auxiliary gyroscope control signal $\boldsymbol{\omega}_k$ is given as follows. The control can be embedded into the time-varying dynamical model such that $\mathbf{f}_k(\mathbf{x}, \boldsymbol{\varepsilon}_k) := \mathbf{f}(\mathbf{x}, \boldsymbol{\omega}_k, \boldsymbol{\varepsilon}_k)$. The state-space model takes the canonical form:

$$\begin{aligned} \mathbf{x}_k &= \mathbf{f}_k(\mathbf{x}_{k-1}, \boldsymbol{\varepsilon}_k), \\ \mathbf{y}_k &= \mathbf{h}_k(\mathbf{x}_k) + \boldsymbol{\gamma}_k, \end{aligned} \tag{6}$$

where $\mathbf{x}_k \in \mathbb{R}^n$ is the state at time step $t_k$, $k = 1, 2, \ldots$, $\mathbf{y}_k \in \mathbb{R}^m$ is a measurement, $\boldsymbol{\varepsilon}_k \sim \mathrm{N}(\mathbf{0}, \mathbf{Q}_k)$ is the Gaussian process noise, and $\boldsymbol{\gamma}_k \sim \mathrm{N}(\mathbf{0}, \boldsymbol{\Sigma}_k)$ is the Gaussian measurement noise. The dynamics and measurements are specified in terms of the dynamical model function $\mathbf{f}(\cdot, \cdot, \cdot)$ and the measurement model function $\mathbf{h}_k(\cdot)$.

In this work we employ the extended Kalman filter (EKF, [2, 13]) which provides a means of approximating the state distributions $p(\mathbf{x} \mid \mathbf{y}_{1:k})$ with Gaussians:

$$p(\mathbf{x}_k \mid \mathbf{y}_{1:k}) \simeq \mathrm{N}(\mathbf{x}_k \mid \mathbf{m}_{k|k}, \mathbf{P}_{k|k}). \tag{7}$$

In the EKF, these approximations are formed by first-order linearizations of the non-linearities. The extended Kalman filtering recursion can be written as follows (see [23] for detailed presentation). The dynamics are incorporated into the *prediction step*:

$$\begin{aligned} \mathbf{m}_{k|k-1} &= \mathbf{f}_k(\mathbf{m}_{k-1|k-1}, \mathbf{0}), \\ \mathbf{P}_{k|k-1} &= \mathbf{F}_\mathbf{x}(\mathbf{m}_{k-1|k-1}) \, \mathbf{P}_{k-1|k-1} \, \mathbf{F}_\mathbf{x}^\mathsf{T}(\mathbf{m}_{k-1|k-1}) + \\ & \quad \mathbf{F}_{\boldsymbol{\varepsilon}}(\mathbf{m}_{k-1|k-1}) \, \mathbf{Q}_k \, \mathbf{F}_{\boldsymbol{\varepsilon}}^\mathsf{T}(\mathbf{m}_{k-1|k-1}), \end{aligned} \tag{8}$$

where the dynamic model is evaluated with the outcome from the previous step and zero noise, and $\mathbf{F}_\mathbf{x}(\cdot)$ denotes the Jacobian matrix of $\mathbf{f}_k(\cdot, \cdot)$ with respect to $\mathbf{x}$ and $\mathbf{F}_{\boldsymbol{\varepsilon}}(\cdot)$ with respect to the process noise $\boldsymbol{\varepsilon}$.

We will also address the special case, where the dynamics are entirely linear, that is $\mathbf{f}_k(\mathbf{x}) = \mathbf{A}_k$. In that case the filter prediction step is entirely given by the standard Kalman filter preduction step:

$$\begin{aligned} \mathbf{m}_{k|k-1} &= \mathbf{A}_k \, \mathbf{m}_{k-1|k-1}, \\ \mathbf{P}_{k|k-1} &= \mathbf{A}_k \, \mathbf{P}_{k-1|k-1} \, \mathbf{A}_k^\mathsf{T} + \mathbf{Q}_k. \end{aligned} \tag{9}$$

Measurement data providing observations of the system state at given time steps are combined with the model in the *update step*:

$$\begin{aligned} \mathbf{v}_k &= \mathbf{y}_k - \mathbf{h}_k(\mathbf{m}_{k|k-1}), \\ \mathbf{S}_k &= \mathbf{H}_\mathbf{x}(\mathbf{m}_{k|k-1}) \, \mathbf{P}_{k|k-1} \, \mathbf{H}_\mathbf{x}^\mathsf{T}(\mathbf{m}_{k|k-1}) + \boldsymbol{\Sigma}_k, \\ \mathbf{K}_k &= \mathbf{P}_{k|k-1} \, \mathbf{H}_\mathbf{x}^\mathsf{T}(\mathbf{m}_{k|k-1}) \, \mathbf{S}_k^{-1}, \\ \mathbf{m}_{k|k} &= \mathbf{m}_{k|k-1} + \mathbf{K}_k \, \mathbf{v}_k, \\ \mathbf{P}_{k|k} &= [\mathbf{I} - \mathbf{K}_k \, \mathbf{H}_\mathbf{x}(\mathbf{m}_{k|k-1})] \, \mathbf{P}_{k|k-1} \\ & \quad [\mathbf{I} - \mathbf{K}_k \, \mathbf{H}_\mathbf{x}(\mathbf{m}_{k|k-1})]^\mathsf{T} + \mathbf{K}_k \, \boldsymbol{\Sigma}_k \, \mathbf{K}_k^\mathsf{T}, \end{aligned} \tag{10}$$

where $\mathbf{H}_\mathbf{x}(\cdot)$ denotes the Jacobian of the measurement model $\mathbf{h}_k(\cdot)$ with respect to the state variables $\mathbf{x}$. The slightly unorthodox form of the last line is known as the Joseph's formula, which both numerically stabilizes updating the covariance and preserves symmetry.

The linearizations inside the extended Kalman filter cause some errors in the estimation. Most notably the estimation scheme does not preserve the norm of the orientation quaternions. Therefore after each update an extra quaternion normalization step is added to the estimation scheme.

## C. Propagation by Gyroscope Prediction

The state holds the internal and external parameters of the camera and the 3D position of the features. The internal parameters and the 3D coordinates of the features are assumed to stay constant (but are still unknown), so their propagation functions are identities. The external parameters (position, orientation) of the camera are not constant and are treated differently.

The position of the camera is modeled as a Wiener velocity process, a commonly used model in tracking and control literature (see, *e.g.*, [23] for details). In order to keep track of the full inertial state the estimate contains both the position and velocity vectors $\mathbf{x} = \begin{pmatrix} \mathbf{p} & \mathbf{v} \end{pmatrix}^\mathsf{T}$ and the acceleration is modeled as a white noise process $\frac{\mathrm{d}^2 x(t)}{\mathrm{d}t^2} = w(t)$, or in state space form as a linear time-invariant stochastic differential equation (independently for each spatial dimension)

$$\frac{\mathrm{d}\mathbf{x}(t)}{\mathrm{d}t} = \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix} \mathbf{x}(t) + \begin{pmatrix} 0 \\ 1 \end{pmatrix} w(t), \tag{11}$$

where $w(t)$ is a realization of the white noise process. In discrete time the system is

$$\mathbf{x}_{k+1} = \mathbf{A}\mathbf{x}_k + \boldsymbol{\varepsilon}_k \tag{12}$$

where $\boldsymbol{\varepsilon}_k \sim \mathrm{N}(\mathbf{0}, \mathbf{Q}_t)$, and $\mathbf{x}_k := \mathbf{x}(t_k)$.

The orientation of the camera is encoded in a unit quaternion, unit quaternions are a direct representation of an axis-angle rotation and can be converted into a rotation matrix using Rodriguez formula.

Given a unit quaternion $\mathbf{q}$ that represents a rotation and a known angular rate $\boldsymbol{\omega}$ in the same frame of reference, its derivative can be expressed as

$$\frac{\mathrm{d}\mathbf{q}(t)}{\mathrm{d}t} = \frac{1}{2}\,\Omega(\boldsymbol{\omega})\,\mathbf{q}(t), \tag{13}$$

where

$$\Omega(\boldsymbol{\omega}) = \begin{pmatrix} 0 & -\boldsymbol{\omega}^{\mathsf{T}} \\ \boldsymbol{\omega} & [\boldsymbol{\omega}]_{\times} \end{pmatrix}. \tag{14}$$

The notation $[w]_{\times}$ is the $3 \times 3$ cross-product matrix (for further details on quaternion modeling, see, *e.g.*, [17]).

Assuming constant rotation rate (during $\Delta t$), the discrete-time system is

$$\mathbf{q}_{k+1} = \exp\left(\frac{\Delta t_k\,\Omega(\boldsymbol{\omega}_k)}{2}\right)\mathbf{q}_k. \tag{15}$$

Since the gyroscope produces rotational rate measurements with known accuracy, it can be propagated into an uncertainty for the quaternion. The gyroscope data is used directly in the prediction as a control signal in a linear Kalman filter.

Putting it all together, the state dynamics are described by

$$\mathbf{f}_k(\mathbf{x}_k, \boldsymbol{\varepsilon}_k) = \mathbf{A}_k\mathbf{x} + \boldsymbol{\varepsilon}_k, \tag{16}$$

where the linear dynamics are given by

$$\mathbf{A}_k = \begin{pmatrix} \mathbf{I}_3 & & & & \\ & \mathbf{I}_3 & \mathbf{I}_3\Delta t_k & & \\ & & \mathbf{I}_3 & & \\ & & & \exp(\frac{\Delta t_k}{2}\,\Omega(\boldsymbol{\omega}_k - \boldsymbol{\omega}_b)) & \\ & & & & \mathbf{I}_3 \end{pmatrix} \tag{17}$$

and

$$\boldsymbol{\varepsilon}_k \sim \mathrm{N}(\mathbf{0}, \mathbf{Q}), \quad \mathbf{Q} = \mathrm{blkdiag}\begin{pmatrix} \mathbf{0}_3 & \mathbf{Q}_t & \mathbf{Q}_q & \mathbf{0}_3 \end{pmatrix}, \tag{18}$$

where $\boldsymbol{\omega}_b$ is the gyroscope bias (estimated off-line) and the $\mathbf{Q}_t$ and $\mathbf{Q}_q$ matrices model the process noise of the translation and rotation, respectively.

The process noise of the translation is derived from the wiener velocity model described above, the process noise of the rotation is propagated from the rotational rate.

### D. Visual Update

The visual update couples all of the state variables. In Figure 2 the visual update occurs every time a new frame has been acquired. The frame is first processed by the feature tracker and inlier detection, and then the feature pixel coordinates are passed to the visual update model which processes an extended Kalman filter update step.

The features are chosen by the Shi–Tomasi *Good features to track* method [25] which determines strong corners in the image. These features are tracked across frames by a pyramidal *Lucas–Kanade tracker* [3, 20]. The *Seven-point algorithm* [9] is used for inlier detection. If a feature is recognized as an inlier, the visual update directly proceeds for the feature. If the feature is not an inlier, the feature is replaced in the state by re-initializing its current position $\mathbf{z}^{(i)}$ estimate (both in terms of state mean and covariance) to uninformative *a priori* values.

| Variable | Initial | Batch opt. | Jia and Evans [15] | Proposed |
|---|---|---|---|---|
| $f_x$ (px) | 75 | 0.05 | 15.19 | 0.36 |
| $f_y$ (px) | 75 | 0.05 | 15.19 | 0.38 |
| $c_x$ (px) | 0.5 | 0.02 | 1.91 | 0.27 |
| $c_y$ (px) | 0.5 | 0.10 | 2.09 | 0.34 |
| $k_1$ | 0.01 | 0.0001 | — | 0.0001 |
| $k_2$ | 0.01 | 0.0095 | — | 0.0095 |

The measurement model function $\mathbf{h}(\cdot)$ is a function of all state variables. The measurement vector $\mathbf{y} \in \mathbb{R}^{2d}$ contains the pixel coordinates, and feature-wise the model can be written as

$$\mathbf{y}^{(i)} = \mathbf{h}^{(i)}(\mathbf{x}) = \mathrm{distort}\left[\mathbf{K}(\mathbf{c})\,\mathbf{E}(\mathbf{p}, \mathbf{q})\begin{pmatrix} \mathbf{z}^{(i)} \\ 1 \end{pmatrix}, \mathbf{c}\right], \tag{19}$$

for $i = 1, 2, \ldots, d$, where $d$ is the number of features being tracked. The measurement model is fully determined by the camera model in Eqs. (1)–(4). For the EKF update, the Jacobian matrix of Eq. (19) must be formed as well. The measurement noise is set to $\boldsymbol{\Sigma} = (2.5\,\mathrm{px})^2\,\mathbf{I}$.

### III. EXPERIMENTS

We demonstrate our self-calibration method both in a simulation setup with known ground-truth values, and empirically on real data. All the experiments runs were performed on the data off-line using a Mathworks MATLAB implementation[1], and benchmarked against the publicly available MATLAB implementation by Jia and Evans [15]. Furthermore, we compared our approach to a bundle-adjustment setup, which was used for checking how well the other methods actually were able to perform.

### A. Simulation Study

A simulation similar to the one described in [15] was performed to compare the results. To evaluate the methods, 100 camera movement paths around a three-dimensional regular point structure were simulated in a Monte Carlo setup. The point coordinates were projected into a virtual camera that followed the paths and continuously turned to look at the point structure. The ground-truth camera model was a zero-skew square pixel camera with a focal length of 575 pixels, an image size of $480 \times 640$ pixels, with the origin at the center of the frame. The distortion parameters were set to be zero in the simulation.

The simulation setup produced tracks for the visual features and gyroscope readings, at 10 Hz and 100 Hz, respectively. The initial state estimates corresponding to the camera parameters were set as follows. The initial guess for the origin of the image plane and focal length were the image center coordinates and 700 px, respectively.

Table I shows the average RMSE results over the 100 simulations for the initial state and three methods run on
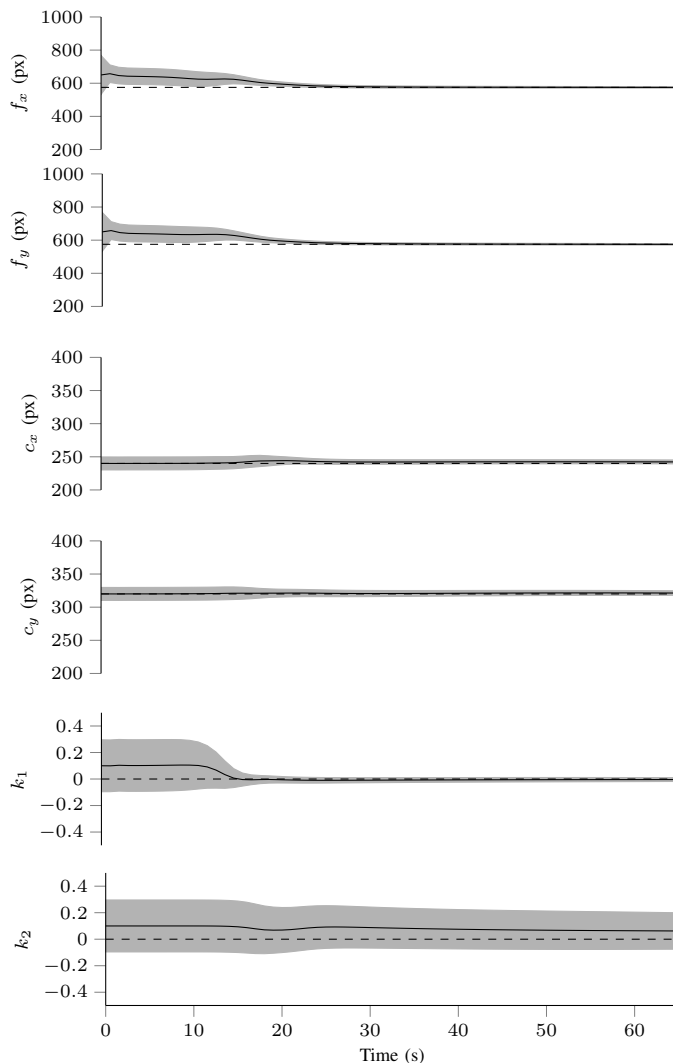
Fig. 3. Evolution of the camera calibration when started from distant initial values in one Monte Carlo simulation. Initial convergence is reached immediately after the first image pair has been observed, and final convergence once the camera has moved sufficiently. The shaded regions depict the 95% uncertainty interval given by the state-space estimation.

the results show that the data is informative about the parameter values, as the optimization converges to the ground-truth values.

The third method was the method by Jia and Evans [15]. We used their autocalibration toolbox[2] implementation for solving the self-calibration problem. As can be seen in the results, there were issues in running these results. As the method by Jia and Evans assumes there to be a lot of features available, we re-ran the results by including almost ten time more features in the simulation. The method keeps oscillating around the correct camera parameter values, but fails to converge. We suspect that the issues are partly related to the use case not being optimal for their method, and partly because of issues in their implementation of the method.

Figure 3 shows the evolution of the camera calibration, when started from distant initial values. The shaded regions depict the 95% uncertainty interval given by the state-space estimation. The model reaches initial convergence immediately after the first image pair has been observed, and final convergence once the camera has moved sufficiently. Note that $k_2$ does not show the convergence of the other parameters. The main reason is that for this case the features are around the center of the image at the beginning, where the $r$ from 4 is small and thus the gradient with respect to $k_2$ is small. For the synthetic case, the distortion parameters are an over-model.

The results show that the proposed method has an RMSE less than one in pixel units for the intrinsic parameters, not far from the results given brute-force calculated bundle-adjustment based calibration results. The method by Jia and Evans is clearly off more, but still did not diverge.

### B. Empirical Tests

In order to demonstrate the proposed online calibration approach on empirical data, we acquired test data in various situations using an Apple iPad Pro 12.9-inch model. Hardware-wise the iPad can be seen as a representative example of a modern-day smart-device. Similar sensors are available in most Android and iOS smartphones and tablets.

The data acquisition was conducted using a custom data capture app implemented in Objective-C. The capture tool app stored the three-axis gyroscope together with associated timestamps to a file in the device. The gyroscope sampling rate was set to 100 Hz. The gyroscope was pre-calibrated before

[2]Online Camera-Gyroscope Auto-Calibration for Cellphones: http://users.ece.utexas.edu/~bevans/papers/2015/autocalibration/

the simulated data. The first column shows the initial RMSE from the initialized parameters. We compare the calibration parameters acquired from three models. The first is the model proposed, which was initialized and run as described earlier in the Methods section.

To analyze the observability of the parameters in the simulated data, we provide results also for a bundle-adjusted batch solution ('Batch opt.' in Table I), which acts as a brute-force baseline for the methods. For this baseline the camera parameters and motion track estimated by our proposed method were used for initializing the bundle adjustment (non-linear minimization of the reprojection error). The bundle adjustment problem was solved by constrained minimization using the MATLAB fmincon Interior Point Algorithm. This bundle adjustment approach was computationally very intensive, but

(a) Cards—A dark indoor data set



(b) Cups—A feature-poor data set
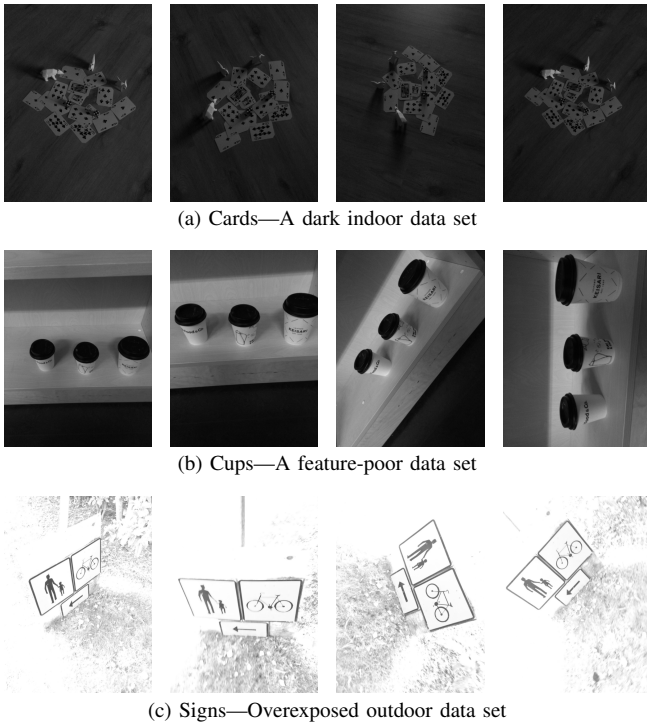


(c) Signs—Overexposed outdoor data set

Fig. 4. Example frames from the empirical tests: (a) is underexposed and with uneven concentration of features, (b) is from a feature-poor indoor scene, and (c) is an overexposed outdoor scene.

the data acquisition for estimating the additive gyroscope bias $\boldsymbol{\omega}_b$.

Simultaneously to the gyroscope capture, device camera frames were read time-locked to the gyroscope observations. The experiments use the rear-facing camera with a resolution of $480 \times 640$ (portrait orientation), grayscale images, exposure time $1/60$, fixed aperture, sensitivity (ISO value) 125, and locked focus at infinity. The camera refresh rate was 10 fps (Hz). The camera frames were stored as H.264 packed video sequences on the device, with exact frame timestamps stored separately for use in the offline run.

Furthermore, camera images of the canonical OpenCV checkerboard pattern were captured for conducting batch calibration of the tablet camera. This calibration was only used for obtaining ground-truth camera parameters to compare against.

For obtaining a set of verstaile use cases, we recorded short sequences of a few controlled static scenes. In these data sets the motion of the camera was smooth in the sense of avoiding hard stops, in order to comply with the constrains of the model. Figure 4 shows example frames from the empirical data sets. They include (a) ill-lit (underexposed) indoor scenes (we name this data set 'cards'), (b) visual feature poor scenes (typical in office environments), and (c) overexposed outdoor scenes.

The evolution of the proposed calibration method is shown in Figure 5, where the shaded area represents the 95% uncertainty interval. The corresponding calibration result is shown in Table II together with the initial values and the checkerboard calibrated ground-truth for the device camera.
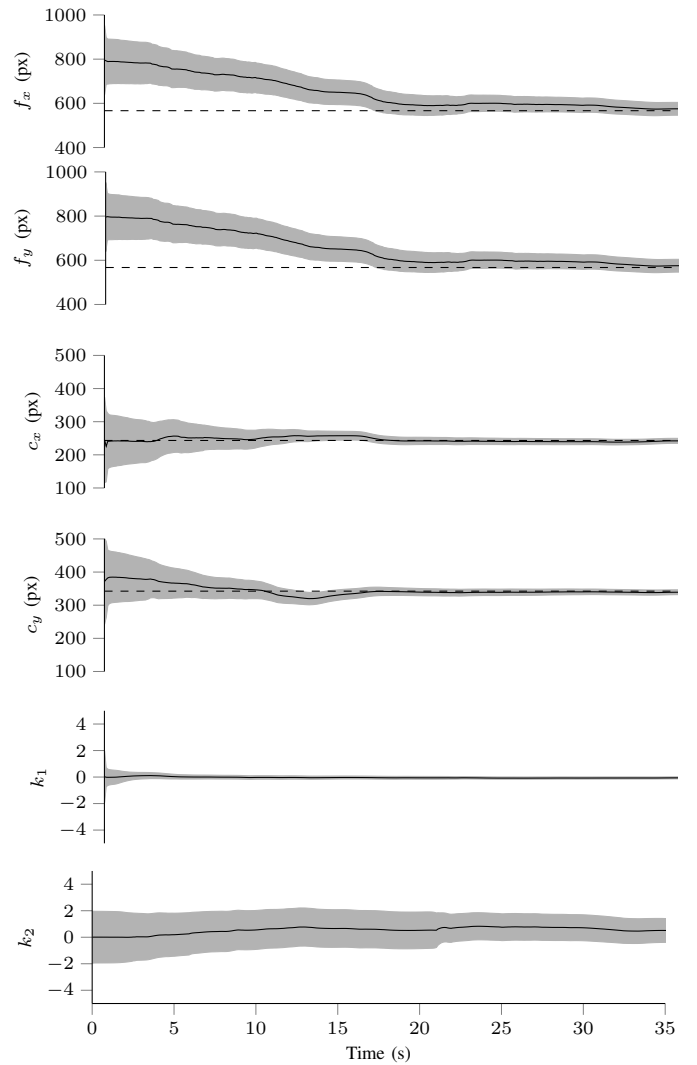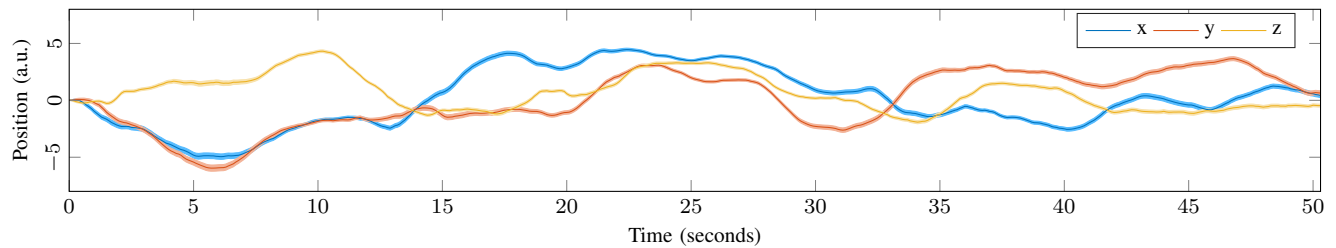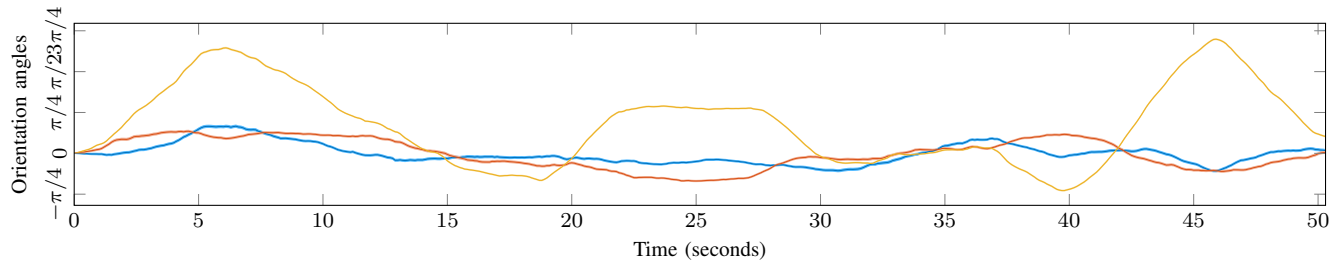


Fig. 5. Evolution of the camera calibration in the 'Cards' data set, when started from distant initial values. The shaded regions depict the 95% uncertainty interval given by the state-space estimation. The convergence occurs once sufficient excitation movement has occurred.

For this run, the evolution of the position estimates and the orientation (transformed into Euler angles) is visualized in Figure 6. Subfigure (a) shows the latent (unobserved) camera position track estimates (the shaded area represents the 95% credible interval) over the time-span of the data. The absolute scale of the movement remains unobserved, which is due to the gyroscope only observing the rotation rate and the camera data being agnostic to the true distance of any feature movement. Subfigure (b) shows the corresponding orientation states, which coincide more clearly with the behavior in Figure 5. The linear parameters appear to reach the right regime after the camera has been rotated sufficiently (*i.e.* the data features sufficient excitation). The non-linear distortion parameters take longer to stabilize.
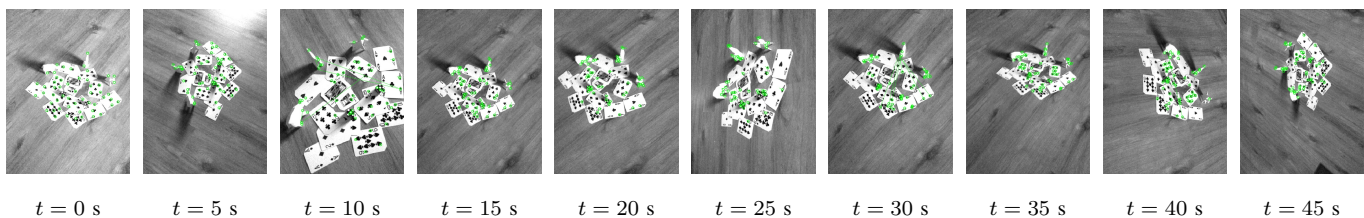
Subfigure 6(c) shows the corresponding input frames in the cards data set with challenging lighting conditions. The green

(a) Latent camera position state estimates



(b) Camera orientation estimates



| $t = 0$ s | $t = 5$ s | $t = 10$ s | $t = 15$ s | $t = 20$ s | $t = 25$ s | $t = 30$ s | $t = 35$ s | $t = 40$ s | $t = 45$ s |

(c) Associated frames at different time points

Fig. 6. State estimates and example input frames for the parameter estimate track results in Fig. 5. Subfigure (a) shows the latent (unobserved) camera position track estimates (the shaded area represents the 95% credible interval) over the time-span of the data. Note that the absolute scale of the movement remains unobserved. (b) The corresponding orientation states, which coincide more clearly with the behavior in Fig. 5. (c) Example input frames in the cards data set with green markers showing the tracked feature positions. Observation times correspond to the tick marks in the above plots.

markers show the current feature point locations in the frames, and their 'tails' (green line) show the point movement since the previous frame (frames sampled at 10 Hz).

For the intrinsic parameters, the estimated parameter values converge within a few pixels to the checkerboard-calibrated ground-truth values. For the non-linear radial distortion parameters the values are also similar. In case of the distortion parameters, the identification might suffer from the low feature concentration towards the edges of the visual frame data.

## IV. DISCUSSION AND CONCLUSION

In this paper we have proposed a method for estimating the intrinsic parameters and lens distortion coefficients for camera calibration. This paper proposed a model for estimating the parameters on the fly by fusing gyroscope and camera data, where the model is based on joint estimation of visual feature positions, camera parameters, and the camera pose, the rotation of which is assumed to follow the movement predicted by the gyroscope.

The estimation procedure is lightweight and performs online using an extended Kalman filter. The strengths of the method is in robustness to feature-poor visual environments and insensitivity to initialization, the aspects which were also

demonstrated in the experiments on simulated and empirical data.

The experiments showed that the method performs well against the current state-of-the-art in gyroscope-aided self-calibration. The results showed that after sufficient motion the convergence is both quick, and it is easy to capture the moment of convergence by monitoring the state variance estimates for the camera parameters.

The empirical tests demonstrated the method using data captured using an Apple iPad Pro. The empirical data the parameter estimates converged rapidly after the system had experienced sufficient motion required for jointly estimating both the camera poses and feature world coordinates.

Figure 4 shows three distinctive data sets; one is underexposed and with uneven concentration of features, one is from a feature-poor indoor scene, and the third is from an overexposed outdoor scene. The proposed method is not sensitive to feature-poor environments as it can rely on a relatively low number of features being tracked—in the simulations only 27 features were used.

This is a clear difference to previous methods, which have been mostly inspired by building on purely machine vision methods. Requiring hundreds of tracked features to

be available works well in controlled environments, in good lighting conditions, and using high-quality camera hardware. Our method focuses on the opposite—low quality data and a low number of features suffices. On the other hand, this comes with a requirement of the features remaining visible for a sufficiently long time period.

The method could be extended to include further parameters. Natural directions of extension would be to also estimate the additive three-axis gyroscope bias, or include further camera parameters such as rolling-shutter timing parameters.

As the method jointly estimates both camera poses, parameters, and feature positions, it can be sensitive to certain use cases and environments. For example, in Figure 4(b) the cups are round and their corresponding features are often lost and picked up again during movement. Short feature tracks inject uncertainty into the estimation scheme in the proposed method, which slows down convergence and misleads the estimation. This could be improved by tuning the method parameters, or introducing visual loop-closures which would reuse the same features when they appear again.

Code implementing the method is available online: https://aaltovision.github.io/camera-gyro-calibration/

### REFERENCES

[1] S. Álvarez, D. Llorca, and M. Sotelo. Hierarchical camera auto-calibration for traffic surveillance systems. *Expert Systems with Applications*, 41(4):1532–1542, 2014.

[2] Y. Bar-Shalom, X.-R. Li, and T. Kirubarajan. *Estimation with Applications to Tracking and Navigation*. Wiley-Interscience, New York, 2001.

[3] J.-Y. Bouguet. Pyramidal implementation of the affine Lucas Kanade feature tracker: Description of the algorithm. Technical report, Intel Corporation, 2001.

[4] J. Civera, D. R. Bueno, A. J. Davison, and J. Montiel. Camera self-calibration for sequential Bayesian structure from motion. In *International Conference on Robotics and Automation (ICRA)*, pages 403–408. IEEE, 2009.

[5] F. Faion, P. Ruoff, A. Zea, and U. D. Hanebeck. Recursive bayesian calibration of depth sensors with non-overlapping views. In *Proceedings of the 15th International Conference on Information Fusion (FUSION)*, pages 757–762, July 2012.

[6] O. D. Faugeras, Q.-T. Luong, and S. J. Maybank. Camera self-calibration: Theory and experiments. In *European Conference on Computer Vision (ECCV)*, pages 321–334, 1992.

[7] M. Goldshtein, Y. Oshman, and T. Efrati. Seeker gyro calibration via model-based fusion of visual and inertial data. In *Proceedings of the 10th International Conference on Information Fusion (FUSION)*, pages 1–8, July 2007.

[8] N. S. Gopaul, J. Wang, and B. Hu. Camera auto-calibration in GPS/INS/stereo camera integrated kinematic positioning and navigation system. *Journal of Global Positioning Systems*, 14(1):3, 2016.

[9] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, second edition, 2004.

[10] A. Heidari, I. Alaei-Novin, and P. Aarabi. Fusion of spatial and visual information for object tracking on iPhone. In *Proceedings of the 16th International Conference on Information Fusion (FUSION)*, pages 630–637, July 2013.

[11] D. Herrera, C. J. Kannala, and J. Heikkila. Forget the checkerboard: Practical self-calibration using a planar scene. In *Winter Conference on Applications of Computer Vision (WACV)*, pages 1–9, 2016.

[12] M. Hwangbo, J.-S. Kim, and T. Kanade. Gyro-aided feature tracking for a moving camera: Fusion, auto-calibration and GPU implementation. *The International Journal of Robotics Research*, 30(14):1755–1774, 2011.

[13] A. H. Jazwinski. *Stochastic Processes and Filtering Theory*. Academic Press, New York, 1970.

[14] C. Jia and B. L. Evans. Online calibration and synchronization of cellphone camera and gyroscope. In *Proceedings of the Global Conference on Signal and Information Processing (GlobalSIP)*, pages 731–734. IEEE, 2013.

[15] C. Jia and B. L. Evans. Online camera-gyroscope autocalibration for cell phones. *IEEE Transactions on Image Processing*, 23(12):5070–5081, 2014.

[16] A. Karpenko, D. Jacobs, J. Baek, and M. Levoy. Digital video stabilization and rolling shutter correction using gyroscopes. Technical Report 2011-3, Stanford University, 2011.

[17] M. Kok. *Probabilistic Modeling for Positioning Applications Using Inertial Sensors*. PhD thesis, Linköping University, Linköping, Sweden, 2014.

[18] F. Li, H. Sekkati, J. Deglint, C. Scharfenberger, M. Lamm, D. Clausi, J. Zelek, and A. Wong. Simultaneous projector-camera self-calibration for three-dimensional reconstruction and projection mapping. *IEEE Transactions on Computational Imaging*, 3(1):74–83, March 2017.

[19] M. Li, H. Yu, X. Zheng, and A. I. Mourikis. High-fidelity sensor modeling and self-calibration in vision-aided inertial navigation. In *International Conference on Robotics and Automation (ICRA)*, pages 409–416. IEEE, 2014.

[20] B. D. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *International Conference on Artificial Intelligence (IJCAI)*, pages 674–679. Vancouver, BC, Canada, 1981.

[21] J. Maye, P. Furgale, and R. Siegwart. Self-supervised calibration for robotic systems. In *IEEE Intelligent Vehicles Symposium (IV)*, pages 473–480, 2013.

[22] H. Ovrén and P.-E. Forssén. Gyroscope-based video stabilisation with auto-calibration. In *International Conference on Robotics and Automation (ICRA)*, pages 2090–2097. IEEE, 2015.

[23] S. Särkkä. *Bayesian Filtering and Smoothing*. Cambridge University Press, 2013.

[24] M. A. Shelley. Monocular visual inertial odometry on a mobile device. Master's thesis, Technical University of Munich, Germany, 2014.

[25] J. Shi and C. Tomasi. Good features to track. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 593–600, 1994.

[26] A. Solin, S. Cortes, E. Rahtu, and J. Kannala. Inertial odometry on handheld smartphones. In *Proceedings of the International Conference on Information Fusion (FUSION)*, 2018.

[27] A. Solin, S. Cortes, E. Rahtu, and J. Kannala. PIVO: Probabilistic inertial-visual odometry for occlusion-robust navigation. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 616–625, 2018.

[28] J. Sun, P. Wang, Z. Qin, and H. Qiao. Effective self-calibration for camera parameters and hand-eye geometry based on two feature points motions. *IEEE/CAA Journal of Automatica Sinica*, 4(2):370–380, 2017.

[29] S. Willi and A. Grundhöfer. Robust geometric self-calibration of generic multi-projector camera systems. In *International Symposium on Mixed and Augmented Reality (ISMAR)*, pages 42–51, 2017.