# Spatiotemporal Calibration of 3D Millimetre-Wavelength Radar-Camera Pairs

Emmett Wise, Qilong Cheng, and Jonathan Kelly

arXiv:2211.01871v4 [cs.RO] 10 Dec 2023

*Abstract*—**Autonomous vehicles (AVs) often depend on multiple sensors and sensing modalities to impart a measure of robustness when operating in adverse conditions. Radars and cameras are popular choices for use in combination; although radar measurements are sparse in comparison to camera images, radar scans are able to penetrate fog, rain, and snow. Data from both sensors are typically fused prior to use in downstream perception tasks. However, accurate sensor fusion depends upon knowledge of the spatial transform between the sensors and any temporal misalignment that exists in their measurement times. During the life cycle of an AV, these calibration parameters may change, so the ability to perform in-situ spatiotemporal calibration is essential to ensure reliable long-term operation. State-of-the-art 3D radar-camera spatiotemporal calibration algorithms require bespoke calibration targets that are not readily available in the field. In this paper, we describe an algorithm for *targetless* spatiotemporal calibration that is able to operate without specialized infrastructure. Our approach leverages the ability of the radar unit to measure its own ego-velocity relative to a fixed, external reference frame. We analyze the identifiability of the spatiotemporal calibration problem and determine the motions necessary for calibration. Through a series of simulation studies, we characterize the sensitivity of our algorithm to measurement noise. Finally, we demonstrate accurate calibration for three real-world systems, including a handheld sensor rig and a vehicle-mounted sensor array. Our results show that we are able to match the performance of an existing, target-based method, while calibrating in arbitrary, infrastructure-free environments.**

*Index Terms*—**Calibration & Identification, Sensor Fusion, Robot Sensing Systems, Radar, Computer Vision**
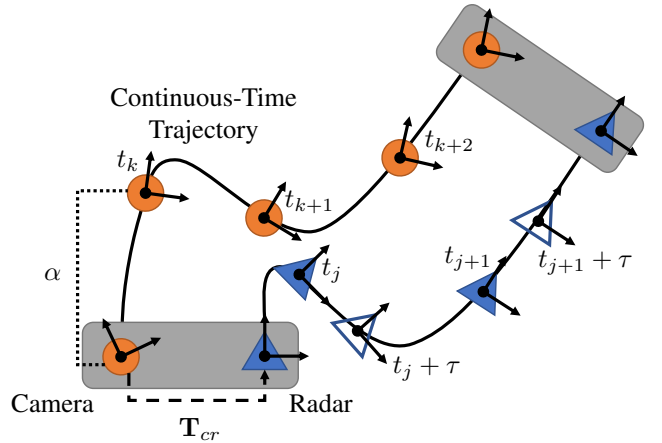


Fig. 1: The radar (triangle) and camera (circle) are assumed to be rigidly connected. Our calibration problem involves estimating the transform between the camera and radar, $\mathbf{T}_{cr}$, the translation scale factor, $\alpha$, for the camera pose measurements, and the temporal offset, $\tau$. The unfilled triangles represent radar measurements at "shifted" points in time due to the offset bias, which must be considered to ensure the correct radar ego-velocity estimate.

## I. INTRODUCTION

The widespread deployment of autonomous vehicles (AVs) depends critically on their ability to operate safely under a range of challenging environmental conditions. To ensure sufficient redundancy, most AV perception systems incorporate multiple sensors and sensing modalities. In this paper, we consider 3D mm-wavelength radar as a complementary sensor to standard cameras for safe AV perception.

The operating principle of mm-wavelength radars (i.e., the active emission of mm-wavelength electromagnetic (EM) radiation) makes these sensors relatively immune to adverse conditions that negatively affect cameras. Radars also provide information that cameras do not, including *range-rate*

measurements of the relative velocity of targets in the environment. However, radar data are much lower resolution and significantly more noisy, than visual measurements under nominal conditions. Together, radars and cameras are highly complementary, providing situational awareness under both nominal and visually-degraded conditions.

To be used jointly in the AV perception stack, radar and camera sensors must be calibrated with respect to each other. The spatial (6-DoF) transform between a radar-camera pair must be known accurately in order to express their data in a common reference frame. An AV may undergo calibration 'at the factory' prior to operation, but maintenance and general wear and tear can alter the spatial calibration parameters. Further performance gains are achieved when the sensor data streams are temporally aligned in addition to spatial alignment. Even when the sensors are externally triggered, internal signal processing delays can result in shifted measurement timestamps. If this time offset is not accounted for, then, for example, moving targets will be *spatially* shifted in the radar and camera data. Further, in some systems, power cycling or reconfiguring the sensors may change the time offset. In turn, temporal calibration may need to be performed routinely to ensure the accuracy and integrity of fused sensor measurements.

An in-situ method to estimate the spatial transformation

Emmett Wise, Qilong Cheng, and Jonathan Kelly are with the Space and Terrestrial Autonomous Robotic Systems Laboratory, University of Toronto, Institute for Aerospace Studies, Toronto, Canada. {<first name>.<last name>@robotics.utias.utoronto.ca}

from the radar to the camera and the temporal offset of the sensor data streams, would enable long-term AV operation in the field. However, existing radar-camera spatial and spatiotemporal calibration algorithms are restricted to certain environments and sensor configurations [1], [2]. Primarily, these methods rely on the assumption that the radar measures 'point-like' reflections from objects. In general, a radar measurement, determined from a reflected EM pulse, is a complex function of the shape, relative orientation, size, and composition of an object [3]. Additionally, multipath reflections introduce outlier measurements of ghost 'objects' [3]. To avoid these problems, specialized trihedral retroreflective radar targets are used to produce the desired point-like radar returns. A visual fiducial can be placed over or alongside a trihedral target, allowing radar-camera measurement correspondences to be established. The use of targets, however, means that calibration must be carried out in specialized areas or with infrastructure that is not usually available during regular AV operation. Additionally, these algorithms require that the radar-camera pair(s) share overlapping fields of view, which may not be possible for all radar-camera systems.

Herein, we extend the method in [4] to jointly estimate the extrinsic calibration parameters and temporal offset of a 3D mm-wavelength radar-camera pair in a fully targetless manner. Importantly, our approach does not require the sensors to share overlapping fields of view. Instead, we use radar measurements to estimate the instantaneous radar ego-velocity — that is, the velocity of the radar unit relative to an external reference frame expressed in the radar reference frame [5]. By relying on velocity information, we remove the need for specialized calibration targets while also avoiding the difficult problems of radar and cross-modal data association. We make the following contributions:

- we extend the work in [4] to enable full spatiotemporal calibration of monocular camera-3D radar pairs in arbitrary configurations;
- we prove that the calibration problem is identifiable and determine the motions that are required for reliable calibration;
- we analyze the accuracy of spatiotemporal calibration with varying amounts of sensor noise through an extensive series of simulation studies; and
- we carry out three different real-world experiments, which demonstrate that our algorithm is able to match the accuracy of an existing, target-based method and that we are able to perform calibration in different environments, including for sensors on board an AV.

In Section II, we survey existing extrinsic and spatiotemporal calibration algorithms for mm-wavelength radar sensors. Section III formulates spatiotemporal calibration as a batch, continuous-time estimation problem. We examine the identifiability of the calibration problem in Section IV. In Section V, we describe two simulation experiments designed to evaluate the robustness of our algorithm. In Section VI, we demonstrate the accuracy and flexibility of our algorithm by reporting on three real-world experiments in different environments. Finally, we summarize our work in Section VII.

## II. RELATED WORK

In this section, we survey spatial and spatiotemporal calibration algorithms that can be applied to mm-wavelength radars in relation to another, complementary sensor. Section II-A reviews algorithms for target-based extrinsic calibration, while Section II-B describes algorithms for target-free, or *targetless*, extrinsic calibration. In Section II-C, we discuss prior work on target-based spatiotemporal calibration.

### A. Target-Based Extrinsic Calibration

Early radar extrinsic calibration algorithms, developed prior to the widespread availability of 3D mm-wavelength radar units, were designed to enable 2D radar-camera data fusion. Many of these early extrinsic calibration techniques operate by computing the projective homography that maps points on the horizontal (sensing) radar plane to points on the camera image plane. Because radar sensors are inherently noisy, most calibration algorithms require specialized trihedral reflectors, shown in Fig. 8, that produce coincident, point-like 'signals' in both the radar and camera data, making the correspondence problem easier to solve [6]–[9]. Although 2D radar sensors are not able to properly measure the elevation of remote targets, they do detect targets at a small elevation angle above the radar horizontal plane. Since the distance to off-plane targets will be slightly different, accurate calibration requires that detected reflectors *do* lie on the radar horizontal plane. Sugimoto et al. [6] constrain the trihedral reflector position using the radar return signal strength. During calibration, the approach in [6] filters radar-camera measurement pairs by return intensity; the intensity is maximal for reflectors that lie on the horizontal plane.

More recent 2D radar extrinsic calibration algorithms often minimize a type of 'reprojection error,' that is, the error in the alignment of identifiable objects that appear within both sensors' fields of view. Kim et al. [10] leverage reprojection error to estimate the radar-to-camera transform but assume that radar measurements are strictly constrained to the zero-elevation plane. El Natour et al. [11] determine the radar-to-camera transform by intersecting backprojected camera rays with the 3D 'arcs' along which the 2D radar measurements must lie. Domhof et al. [12] use a specialized calibration target that provides scale for the camera measurement, enabling extrinsic calibration via point cloud alignment. Peršić et al. [13] also perform extrinsic calibration via 3D point cloud alignment but improve overall accuracy by modelling the relationship between target return intensity and elevation angle. The 'homography' and 'reprojection' methods are summarized and compared by Oh et al. in [14], where the authors conclude that both have similar performance. Due to the infrastructure needs (i.e., specialized targets), the methods above are restricted to sensor pairs that share overlapping fields of view. This requirement may be impossible to satisfy for certain sensor configurations. By leveraging constraints induced by the motion of a rigidly-connected radar-camera pair, we are able to calibrate sensors that do not share overlapping fields of view. Further, our approach operates without added infrastructure, enabling calibration under a wider range of conditions.

### B. Target-Free Extrinsic Calibration

Some extrinsic calibration algorithms do not require specialized retroreflective targets. Schöller et al. [15] train a neural network end-to-end to regress a rotation correction from raw camera images and radar data, for example. Peršić et al. [16] estimate the yaw angles (only) between radar, camera, and lidar sensors by aligning the trajectories of tracked objects. Both of these methods require manual measurement of the translation parameters and overlapping fields of view.

Heng [17] presents the first reprojection error-based 3D radar-lidar extrinsic calibration algorithm that does not require specialized targets or overlapping sensor fields of view. The approach in [17] estimates the extrinsic calibration between several lidar units and, using a known vehicle trajectory, constructs a 3D point cloud map. The radar-lidar extrinsic calibration parameters are then determined by minimizing two weighted residuals: i) the distance from the radar point measurements to the closest plane in the lidar map and ii) the radial velocity error. However, this method requires the construction of a dense lidar map and known vehicle poses.

Instead of using feature positions, a subset of extrinsic calibration algorithms fuse ego-velocity and ego-motion measurements from the radar and second sensor, respectively. Since the motion of each sensor is estimated separately, these methods do not perform radar or cross-modal data association and are inherently 'target-free.' Kellner et al. [18] estimate the rotation between a car-mounted 2D radar and an inertial measurement unit (IMU) by minimizing the difference in estimated lateral velocities expressed in the radar frame. While the radar ego-velocity measurements provide lateral velocity directly, determining the lateral velocity of the radar from IMU measurements requires both the IMU angular velocity and accurate knowledge of the radar-IMU translation. Doer et al. [19] extend the approach in [18] to estimate the full extrinsic calibration for a 3D radar-IMU pair. Their method is able to achieve a spatial calibration accuracy of 5 cm and 5° when using simulated, low-noise (our designation, see Section V) radar measurements. Wise et al. [4] perform extrinsic calibration in continuous time using instantaneous radar ego-velocity measurements and camera egomotion measurements. Given an unknown but fixed temporal offset, the spatial calibration parameters estimated by this method are within 3 cm and 1°, per axis, of those determined by [2]. All of these techniques rely on ad-hoc temporal calibration schemes. Herein, we incorporate a principled temporal calibration method.

### C. Target-Based Spatiotemporal Calibration

To date, only two radar spatiotemporal calibration algorithms have appeared in the literature, by Lee et al. [1] and by Peršić et al. [2]. The algorithm in [1] first calibrates the 2D radar-lidar spatial transform using the method of Peršić et al. [13]. As a second step, the lidar measurements are expressed in the radar reference frame, and the azimuth error to distant targets is minimized to determine the temporal offset between the sensor data streams. Peršić et al. [2] represent the trajectory of a target moving through the fields of view of multiple sensors using a continuous-time Gaussian process model.

This representation allows their algorithm to estimate the spatiotemporal calibration parameters by aligning the sensors' trajectories. In general, jointly determining all parameters as part of one maximum likelihood estimation problem yields superior accuracy [2], [20]. Notably, since the methods in [1] and [2] rely on known targets, they have the same limitations as the methods discussed in Section II-A.

## III. METHODOLOGY

We formulate radar-to-camera spatiotemporal calibration as a continuous-time batch estimation problem. In Section III-A, we describe the mathematical notation used throughout the paper. We choose to parameterize the smooth radar-camera trajectories using continuous-time B-splines; we review the properties of this representation in Section III-B. In Section III-C, we derive our radar and camera measurement models. With the necessary preliminaries in place, we then define the full estimation problem in Section III-D.

### A. Notation

Latin and Greek letters (e.g., $a$ and $\alpha$) denote scalar variables, while boldface lower- and uppercase letters (e.g., $\mathbf{x}$ and $\boldsymbol{\Theta}$) denote vectors and matrices, respectively. A parenthesized superscript pair, for example, $\mathbf{A}^{(i,j)}$, indicates the $i$th row and the $j$th column of the matrix $\mathbf{A}$. A three-dimensional reference frame is designated by $\underrightarrow{\mathcal{F}}$. The translation vector from point $a$ (often a reference frame origin) to $b$, expressed in $\underrightarrow{\mathcal{F}}_a$, is denoted by $\mathbf{r}_a^{ba}$. The translational velocity vector of point $b$ relative to point $a$, expressed in $\underrightarrow{\mathcal{F}}_a$, is denoted by $\mathbf{v}_a^{ba}$. The angular velocity of frame $\underrightarrow{\mathcal{F}}_a$ relative to a frame $\underrightarrow{\mathcal{F}}_i$, expressed in $\underrightarrow{\mathcal{F}}_a$, is denoted by $\boldsymbol{\omega}_a^{ai}$.

We denote rotation matrices by $\mathbf{R}$. For example, $\mathbf{R}_{ab} \in \mathrm{SO}(3)$ defines the rotation from $\underrightarrow{\mathcal{F}}_b$ to $\underrightarrow{\mathcal{F}}_a$. We reserve $\mathbf{T}$ for $\mathrm{SE}(3)$ transformation matrices. For example, $\mathbf{T}_{ab}$ is the $4 \times 4$ homogeneous matrix that defines the rigid-body transform from frame $\underrightarrow{\mathcal{F}}_b$ to $\underrightarrow{\mathcal{F}}_a$. Our $\mathrm{SE}(3)$ matrix entries will generally be functions of time; we denote the transform from frame $\underrightarrow{\mathcal{F}}_b$ to $\underrightarrow{\mathcal{F}}_a$ at time $t$ by

$$\mathbf{T}_{ab}(t) = \begin{bmatrix} \mathbf{R}_{ab}(t) & \mathbf{r}_a^{ba}(t) \\ \mathbf{0}^T & 1 \end{bmatrix}, \qquad (1)$$

where $\mathbf{R}_{ab}(t) \in \mathrm{SO}(3)$ and $\mathbf{r}_a^{ba}(t) \in \mathbb{R}^3$. We use $\mathbf{I}_n$ to denote the $n$-by-$n$ identity matrix.

The unary operator $^\wedge$ acts on $\mathbf{r} \in \mathbb{R}^3$ to produce a skew-symmetric matrix such that $\mathbf{r}^\wedge \mathbf{s}$ is equivalent to the cross product $\mathbf{r} \times \mathbf{s}$. The operators $\exp(\cdot)$ and $\log(\cdot)$ map from the Lie algebra $\mathfrak{so}(3)$ to the Lie group $\mathrm{SO}(3)$ and vice versa, respectively [21].

### B. Continuous-Time Trajectory Representation

Temporal calibration is most easily formulated as a continuous-time problem, in part because the batch optimization procedure incrementally time-shifts the measurements from one sensor. In turn, we require the ability to query the pose of the radar or the camera at arbitrary points in time.

To enable this, we parameterize the trajectory of the radar-camera pair using the B-spline representation from Sommer et al. [22]. This representation is briefly reviewed below. We refer readers to Sommer et al. [22], de Boor [23], and Qin [24] for additional details.

A B-spline of order $k$ is a function of one continuous parameter (e.g., time) and a finite set of control points; for brevity, we restrict our example here to control points $\{\mathbf{p}_0, \ldots, \mathbf{p}_N \mid \mathbf{p}_i \in \mathbb{R}^d\}$. In a uniformly-spaced B-spline, each control point is assigned a time (or *knot*) $t_i = t_0 + i\Delta t$, where $t_0$ marks the beginning of the spline and $\Delta t$ is the time between knots. Evaluating a $k^{\text{th}}$ order B-spline at time $t$, where $t_i \leq t < t_{i+1}$, requires the set of $k$ control points over the knot sequence $t_i, \ldots, t_{i+k-1}$. As a result, the end point of a B-spline of length $N$ and order $k$ is at time $t_{N-k+1}$.

The first step in computing the value of a $k^{\text{th}}$ order B-spline at time $t$ is to convert $t$ to the 'normalized' time $u = \frac{t-t_i}{t_{i+1}-t_i}$. Given $u$, the value of the $k^{\text{th}}$ order B-spline is defined as

$$\mathbf{p}(u) = \begin{bmatrix} \mathbf{p}_i & \mathbf{d}_1^i & \ldots & \mathbf{d}_{k-1}^i \end{bmatrix} \tilde{\mathbf{M}}_k \mathbf{u}, \qquad (2)$$

where $\mathbf{u}^T = \begin{bmatrix} 1 & u & u^2 & \ldots & u^{k-1} \end{bmatrix}$ and $\mathbf{d}_j^i = \mathbf{p}_{i+j} - \mathbf{p}_{i+j-1}$. The elements of the $k \times k$ mixing matrix, $\tilde{\mathbf{M}}_k$, are defined by,

$$\tilde{\mathbf{M}}_k^{(a,n)} = \sum_{s=a}^{k-1} m_k^{(s,n)}, \qquad (3)$$

$$m_k^{(s,n)} = \frac{C_{k-1}^n}{(k-1)!} \sum_{l=s}^{k-1} (-1)^{l-s} C_k^{l-s} (k-1-l)^{k-1-n} \qquad (4)$$

$$a, s, n \in \{0, \ldots, k-1\},$$

where scalar $C_j^i = \frac{j!}{i!(j-i)!}$. Substituting $\boldsymbol{\lambda}(u) = \tilde{\mathbf{M}}_k \mathbf{u}$ into Equation (2) results in

$$\mathbf{p}(u) = \mathbf{p}_i + \sum_{j=1}^{k-1} \lambda_j(u) \mathbf{d}_j^i. \qquad (5)$$

Equation (5) can describe the smooth translation of a rigid-body in continuous time (see Figure 10 in Section VI for an example).

While our development above focuses on vector space splines, B-splines can also be defined on Lie groups, including the group SO(3) of rotations,

$$\mathbf{R}(u) = \mathbf{R}_i \prod_{j=1}^{k-1} \exp(\lambda_j(u) \boldsymbol{\phi}_j^i), \qquad (6)$$

where $\mathbf{R}_i$ is a control point of the rotation spline and $\boldsymbol{\phi}_j^i = \log(\mathbf{R}_{i+j-1}^T \mathbf{R}_{i+j})$. We use two B-splines, one on SO(3) and one on $\mathbb{R}^3$, as our complete continuous-time representation of the radar-camera trajectory.

### C. Sensor Measurement Models

In order to perform spatiotemporal calibration, we require a measurement model for the radar unit. Radars emit EM waves that reflect off of surfaces in the environment. Due to the relatively long EM wavelength used by radars, the reflected "location" of a target can vary based on the relative orientation between the radar and target [3]. Additionally,
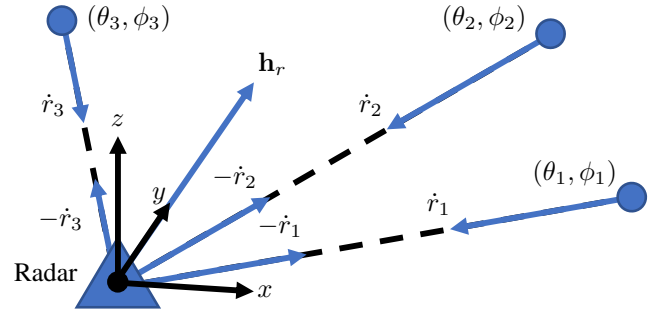


Fig. 2: Illustration of our radar measurement model. The radar EM wave reflects off of three (or more) non-collinear, stationary landmarks in the environment, yielding azimuth, elevation, and range-rate measurements to each landmark. Using these data, we estimate the radar velocity relative to the world reference frame expressed in the radar reference frame.

multipath reflections can occur when the wave bounces off of multiple surfaces before returning to the radar receiver, which can bias measurements of targets and introduce false detections [3].

For each received reflection $l$ from an environmental feature (i.e., an object that we identify as a *landmark*), the radar measures the range $r_l$, azimuth $\theta_l$, elevation $\phi_l$, and range-rate $\dot{r}_l$. We assume that the observed landmarks are stationary with respect to a world frame, $\mathcal{F}_w$; measurements are resolved in the radar reference frame, $\mathcal{F}_r$. The range-rate measurement is the dot product between the velocity of the radar unit itself, $\mathbf{h}_r \in \mathbb{R}^3$, and the unit vector $\hat{\mathbf{r}}_l \in \mathrm{S}^2$ defined by $\theta_l$ and $\phi_l$. Given radar measurements to at least three non-collinear, stationary landmarks, the unit direction vectors and their associated range-rates can be used to reconstruct the radar velocity, $\mathbf{h}_r$, as shown in Figure 2.

Stahoviak [5] and Doer et al. [25] demonstrate that, given $N > 3$ landmarks, one can estimate the ego-velocity of the radar by solving the over-constrained linear least-squares problem

$$\mathbf{h}_r^\star = \min_{\mathbf{h}_r} \ \mathbf{e}_{ego}^T \mathbf{e}_{ego}, \qquad (7)$$

where

$$\mathbf{e}_{ego} = \mathbf{H}\mathbf{x} - \mathbf{y} = \begin{bmatrix} \hat{\mathbf{r}}_0^T \\ \vdots \\ \hat{\mathbf{r}}_N^T \end{bmatrix} \mathbf{h}_r - \begin{bmatrix} \dot{r}_0 \\ \vdots \\ \dot{r}_N \end{bmatrix}. \qquad (8)$$

Equation (8) has its specific form because we wish to estimate the velocity of the radar with respect to the static world frame — not vice versa. The estimated ego-velocity covariance is

$$\boldsymbol{\Sigma}_v = \frac{(\mathbf{e}_{ego}^T \mathbf{e}_{ego})(\mathbf{H}^T \mathbf{H})}{N-3}. \qquad (9)$$

We use RANSAC [5], [25] and radar cross-section thresholding to remove outliers. The two main sources of outliers are targets that move relative to the inertial reference frame and spurious multipath reflections. Empirically, RANSAC successfully eliminates outliers from these two sources if the range-rates to the targets deviate significantly from other stationary landmarks that are close in terms of angle. However, there are two subtle cases where multipath returns may appear to

be valid measurements of stationary landmarks. In the first case, the difference between the transmission and return angles is small. In the second case, the transmission and return angles are symmetric about the radar ego-velocity direction (see Section 8.9 in [3]). The returns from these reflections have a small radar cross-section and are rejected by cross-section thresholding.

Leveraging our continuous-time trajectory representation, the measurement model for the radar ego-velocity at time $t_j$ is

$$\mathbf{h}_{r_j} = -\dot{\mathbf{r}}_r^{wr}(t_j + \tau) - \boldsymbol{\omega}_r^{rw}(t_j + \tau)^{\wedge}\mathbf{r}_r^{wr}(t_j + \tau)$$
$$+ \mathbf{n}_{v_j}, \qquad (10)$$
$$\mathbf{n}_{v_j} \sim \mathcal{N}\left(\mathbf{0}_{3\times1}, \boldsymbol{\Sigma}_{v_j}\right),$$

where $\tau$ is the temporal offset of the radar measurements relative to the camera measurements and $\mathbf{n}_{v_j}$ is the radar velocity measurement noise term. We assume that the noise is a zero-mean Gaussian with covariance matrix $\boldsymbol{\Sigma}_{v_j}$ (see Equation (9)). From the radar ego-velocity model, the error residual is

$$\mathbf{e}_{v_j} = \mathbf{h}_{r_j} + \dot{\mathbf{r}}_r^{wr}(t_j + \tau) +$$
$$\boldsymbol{\omega}_r^{rw}(t_j + \tau)^{\wedge}\mathbf{r}_r^{wr}(t_j + \tau) - \mathbf{n}_{v_j}. \qquad (11)$$

To estimate the ego-motion of the camera, we use a monocular simultaneous localisation and mapping (SLAM) algorithm that operates independently of the radar. By observing fixed landmarks in the environment, monocular SLAM is capable of determining the transformation between the camera reference frame, $\underrightarrow{\mathcal{F}}_c$, and the world frame, $\underrightarrow{\mathcal{F}}_w$, up to an unknown scale factor $\alpha$ [26]. Our (scaled) camera pose measurement model is given by

$$\mathbf{R}_{cw,t_k} = \exp(\mathbf{n}_{r,k})\mathbf{R}_{cr}\,\mathbf{R}_{wr}(t_k)^T, \qquad (12)$$
$$\mathbf{n}_{r,k} \sim \mathcal{N}\left(\mathbf{0}_{3\times1}, \boldsymbol{\Sigma}_r\right),$$

$$\mathbf{r}_{c,t_k}^{wc} = \alpha(\mathbf{R}_{cr}\mathbf{r}_r^{wr}(t_k) + \mathbf{r}_c^{rc}) + \mathbf{n}_{t,k}, \qquad (13)$$
$$\mathbf{n}_{t,k} \sim \mathcal{N}\left(\mathbf{0}_{3\times1}, \boldsymbol{\Sigma}_t\right),$$

where $\mathbf{n}_{r,k}$ and $\mathbf{n}_{t,k}$ are zero-mean Gaussian noise terms of covariances matrice $\boldsymbol{\Sigma}_r$ and $\boldsymbol{\Sigma}_t$ for the camera rotation and translation measurements. The resulting error equations are

$$\mathbf{e}_{r,t_k} = \log(\mathbf{R}_{cw,t_k}\mathbf{R}_{wr}(t_k)\,\mathbf{R}_{cr}^T), \qquad (14)$$
$$\mathbf{e}_{t,t_k} = \mathbf{r}_{c,t_k}^{wc} - \alpha(\mathbf{R}_{cr}\mathbf{r}_r^{wr}(t_k) + \mathbf{r}_c^{rc}) - \mathbf{n}_{t,k}. \qquad (15)$$

We note that a monocular visual odometry (VO) algorithm (i.e., localization without loop closure) can provide camera ego-motion estimates, but visual drift will bias these estimates and decrease calibration accuracy. Also, the measured radar ego-velocity is a local property of a trajectory and cannot fully correct for pose errors induced by visual drift.

### D. The Spatiotemporal Calibration Problem

The set of parameters, $\mathbf{x}$, that we wish to estimate are the spline control points ($\mathbf{r}_{0...N} \in \mathbb{R}^3$, $\mathbf{R}_{0...N} \in \mathrm{SO}(3)$), the extrinsic calibration parameters ($\mathbf{R}_{cr}, \mathbf{r}_c^{rc}$), the camera translation scale factor ($\alpha$), and the temporal offset ($\tau$),

$$\mathbf{x} = \{\mathbf{r}_0, \quad \ldots, \quad \mathbf{r}_N, \quad \mathbf{R}_0, \quad \ldots, \quad \mathbf{R}_N,$$
$$\mathbf{R}_{cr}, \quad \mathbf{r}_c^{rc}, \quad \alpha, \quad \tau\}. \qquad (16)$$

Given $N_r$ radar measurements and $N_c$ camera measurements, we minimize the following cost function,

$$\mathbf{x}^{\star} = \min_{\mathbf{x}} \sum_{j=1}^{N_r} \mathbf{e}_{v_j}^T \boldsymbol{\Sigma}_{v_j}^{-1}\mathbf{e}_{v_j} +$$
$$\sum_{k=1}^{N_c} \mathbf{e}_{r,t_k}^T \boldsymbol{\Sigma}_r^{-1}\mathbf{e}_{r,t_k} + \mathbf{e}_{t,t_k}^T \boldsymbol{\Sigma}_t^{-1}\mathbf{e}_{t,t_k}. \qquad (17)$$

We perform this minimization using the Ceres solver, a standard nonlinear least squares solver [27]. The ability to calibrate all of the relevant parameters depends upon the identifiability of problem, which we discuss in the next section.

## IV. IDENTIFIABILITY

In this section, we show that the calibration problem is identifiable given sufficient excitation of the radar-camera system. Our approach is to determine the observability, or 'instantaneous identifiability,' of the system at several different points in time, assuming that the system follows a varying trajectory. We consider local identifiability (cf. locally weak observability) along a trajectory segment in Section IV-B, after introducing the requisite observability rank condition in Section IV-A. A similar approach has been taken in [28] and [29] and elsewhere. In Section IV-C, we describe several 'degenerate' motions for which the identifiability condition does not hold. We leave the complete characterization of the set of unidentifiable trajectories as future work.

### A. The Observability Rank Condition

We make use of the criterion from Hermann and Krener [30] as part of our identifiability analysis. A system $S$, written in control-affine form as

$$S \begin{cases} \dot{\mathbf{x}} = \mathbf{f}_0(\mathbf{x}) + \sum_{j=1}^{p} \mathbf{f}_j(\mathbf{x})u_j \\ \mathbf{y} = \mathbf{h}(\mathbf{x}) \end{cases}, \qquad (18)$$

with the drift vector field $\mathbf{f}_0(\mathbf{x})$ and control inputs $u_j$ (for $j = 1, \ldots, p$), is locally weakly observable if the matrix $\mathbf{O}$ of the gradients of the Lie derivatives with respect to the system state has full column rank.

The Lie derivative, or directional derivative, of a smooth scalar function $h$ with respect to the smooth vector field $\mathbf{f}$ at the point $\mathbf{x}$ is

$$L_{\mathbf{f}}h(\mathbf{x}) = \nabla_{\mathbf{f}}h(\mathbf{x}) = \frac{\partial h(\mathbf{x})}{\partial \mathbf{x}}\mathbf{f}(\mathbf{x}). \qquad (19)$$

The $n^{th}$ Lie derivative of $h$ with respect to $\mathbf{x}$ along $\mathbf{f}$ is defined recursively as

$$L_{\mathbf{f}}^n h(\mathbf{x}) = \frac{\partial L_{\mathbf{f}}^{n-1} h(\mathbf{x})}{\partial \mathbf{x}}\mathbf{f}(\mathbf{x}), \qquad (20)$$

where $L^0 h(\mathbf{x}) = h(\mathbf{x})$. We note that the matrix $\mathbf{O}$ has an infinite number of rows, but it is sufficient to show that a finite subset of rows yield a matrix of full column rank.

### B. Identifiability of Radar-Camera Calibration

We begin by simplifying the state (and parameter) vector that we aim to estimate. We are able to measure the camera

pose up to scale [26] and the radar velocity in the radar frame [5]. Since we are working in continuous time (or, roughly equivalently, if there are a sufficient number of closely-spaced radar and camera measurements), then the scaled velocity of the camera in the camera reference frame $\alpha \mathbf{v}_c^{cw}(t_i)$, the angular velocity of the camera $\boldsymbol{\omega}_c(t_i)$ in the camera frame, the radar velocity $\mathbf{v}_r^{rw}(t_i + \tau)$ in the radar frame, and the time derivative of the radar velocity $\dot{\mathbf{v}}(t_i + \tau)$ in the radar frame are all available. For the purposes of identifiability, we are able to define the following, modified measurement model,

$$\mathbf{h}(t_i) = \alpha(\mathbf{R}_{cr}\mathbf{v}_r^{rw}(t_i + \tau) - \boldsymbol{\omega}_c(t_i)^{\wedge}\mathbf{r}_c^{rc}), \qquad (21)$$

where $\mathbf{h}(t_i)$ is the scaled linear velocity of the camera ($\mathbf{v}_c^{cw}$) and $\boldsymbol{\omega}_c$ is the angular velocity of the camera, both relative to the camera frame. This modified measurement model does not directly rely on the pose of the radar, thus simplifying the set of parameters that we wish to determine to

$$\tilde{\mathbf{x}} = \{\mathbf{r}_c^{rc}, \quad \mathbf{R}_{cr}, \quad \alpha, \quad \tau\}. \qquad (22)$$

To decrease the notational burden, we drop the superscripts and subscripts defining the velocities and extrinsic transform parameters. The gradient of the zeroth-order Lie derivative of the $i$th measurement is

$$\nabla_{\tilde{\mathbf{x}}} L_0 \mathbf{h}(t_i) = \begin{bmatrix} -\alpha\boldsymbol{\omega}(t_i)^{\wedge} & -\alpha(\mathbf{R}\mathbf{v}(t_i + \tau))^{\wedge}\mathbf{J} \\ \mathbf{R}\mathbf{v}(t_i + \tau) - \boldsymbol{\omega}(t_i)^{\wedge}\mathbf{r} & \alpha\mathbf{R}\dot{\mathbf{v}}(t_i + \tau) \end{bmatrix}, \qquad (23)$$

where $\mathbf{J}$ is the Lie algebra left Jacobian of $\mathbf{R}_{cr}$ [21]. Since the parameters of interest are constant with respect to time, we are able to stack the gradients of several Lie derivatives (at different points in time) to form the observability matrix,

$$\mathbf{O} = \begin{bmatrix} \nabla_{\tilde{\mathbf{x}}} L_0 \mathbf{h}(t_1) \\ \nabla_{\tilde{\mathbf{x}}} L_0 \mathbf{h}(t_2) \\ \nabla_{\tilde{\mathbf{x}}} L_0 \mathbf{h}(t_3) \end{bmatrix}, \qquad (24)$$

which, using block Gaussian elimination, can be shown to have full column rank when three or more sets of measurements are available.[1]

Two comments regarding the analysis are in order. First, we note that the analysis is simplified by considering the modified measurement equation only, which avoids the use of higher-order Lie derivatives. Second, there is a subtlety involved in stacking the gradients of the Lie derivatives at different points in time. The modified measurement equation depends upon the time derivatives of the camera pose and the radar ego-velocity—this implies that, although we do not consider specific control inputs, the system dynamics must be non-null. Stated differently, varied motion of the radar-camera pair is necessary to ensure identifiability; we discuss this requirement further in the next section. Further, it is worth noting that the the observation times must span the temporal offset period [31].

### C. Degenerate Motions

There are motions that cause the matrix $\mathbf{O}$ in Equation (24) to lose full column rank. First, the Lie derivatives include

---

[1]We omit the full derivation for brevity, and note that the rank condition can be verified in this case using any symbolic algebra package.

linear and rotational velocities and accelerations, so the matrix will lose full column rank when the system is stationary with respect to the world frame or moving with constant linear or angular velocity. Second, in Wise et al. [4], we showed that the system must undergo rotation about two nonparallel axes in order for the observability matrix to be full rank. This requirement also applies to the present analysis. To show this, we can align the angular velocity and angular acceleration vectors by substituting $\boldsymbol{\alpha}_c^{ci} = \eta \boldsymbol{\omega}_c^{cw}$, where $\eta$ is an arbitrary constant, into Equation (41). This substitution is equivalent to asserting that the system rotates about one axis only, resulting in an observability matrix that is rank-deficient.

## V. SIMULATION STUDIES

In order to test the robustness of our algorithm to measurement noise, we carried out a series of simulation studies. We generated two simulated camera-radar datasets using two different trajectories and with varying amounts of (synthetic) noise (see Figures 3 and 4). The nominal (noise-free) trajectories were selected to ensure sufficient excitation of the camera-radar pair. The median linear and rotational velocities for the trajectory shown in Figure 3 were, respectively, higher and lower than the velocities for the trajectory shown in Figure 4.

After constructing the trajectories, we computed the simulated radar ego-velocity and camera pose measurements. Since radar measurements are antenna configuration- and environment-specific, these measurements were not generated at the EM propagation level. Radar ego-velocity measurements (i.e., $\mathbf{h}_{r_k}$) were computed using the known linear and rotational velocities defined by the trajectory. Consequently, our simulated radar measurements generalize to any radar and environment that produce an unbiased 3D ego-velocity estimate. Simulated camera pose measurements (i.e., $\mathbf{R}_{cw,t_k}$ and $\mathbf{t}_{c,t_k}^{wc}$) were derived from simulated observations of a series of landmark points, arranged in a 2D grid. This configuration of points matches the configuration of a standard 'checkerboard' camera calibration target. In the targetless setting, we can only estimate the position of the camera up to an unknown scale [26], so the checkerboard tracking algorithm is given an incorrect size for the checkerboard squares.

For each individual simulation, we added zero-mean Gaussian noise to the radar ego-velocity measurements ($\boldsymbol{\Sigma}_{v_j} = \sigma_r^2 \mathbf{I}_{3 \times 3}$) and to the camera measurements of the checkerboard corners on the simulated image plane ($\boldsymbol{\Sigma}_{p_k} = \sigma_c^2 \mathbf{I}_{2 \times 2}$). In our experiments, we adjusted the radar ego-velocity noise standard deviation ($\sigma_r$) between 0.05 m/s and 0.15 m/s. Based on our real-world experiments (see Section VI), we have found that the radar ego-velocity measurement noise is closer to the lower end of this range, unless the environment is sparse and too few valid radar returns are captured. We adjusted the standard deviation of the noise added to the measured checkerboard corner coordinates ($\sigma_c$) between 0.2 and 0.4 pixels; these noise levels are similar to the observed noise in our real-world experiments [4].

The error distributions for the spatial calibration parameter estimates ($\mathbf{R}_{cr}, \mathbf{r}_c^{rc}$), scale factor ($\alpha$), and temporal offset ($\tau$) are shown in Figures 5 and 6, across 100 simulation trials
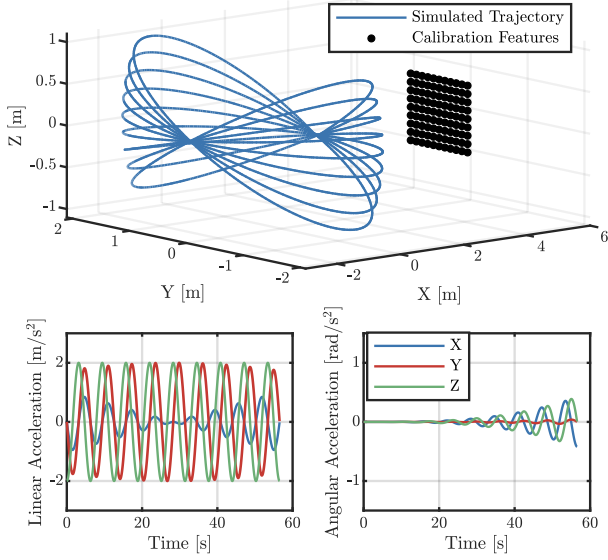
Fig. 3: High linear and low angular velocity trajectory (top) for the simulation experiments, with associated linear (bottom left) and angular acceleration (lower right) plots.
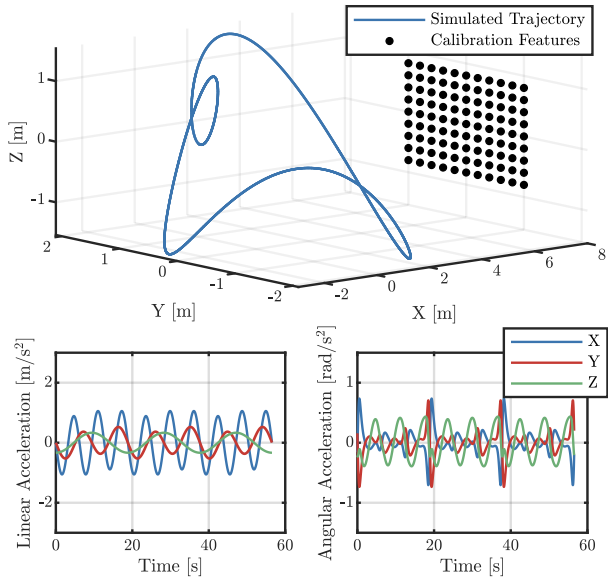


Fig. 4: Low linear and high angular velocity trajectory (top) for the simulation experiments, with associated linear (bottom left) and angular acceleration (lower right) plots.

for each (nominal) trajectory. For the high linear and low angular velocity trajectory, even in the high-noise regime, the error in the rotation and scale estimates remains less than two degrees and one percent, respectively. However, high levels of noise in the radar ego-velocity measurements result in substantially larger (and more widely distributed) errors in the estimate of the relative translation of the sensors and of the temporal offset; the errors can be as large as 15 cm and 30 ms, respectively. This sensitivity indicates that, prior to use in our algorithm, the radar data should be filtered to remove high-noise measurements whenever possible.

If the system follows the low linear and high angular velocity trajectory in Figure 4, then radar data filtering may not

be necessary. As shown in Figure 6, calibrating the radar along this trajectory results in similar scale and rotation estimation accuracy as for the other trajectory, but drastically improves the translation and temporal offset estimates; the errors are within 10 cm and 10 ms, respectively. Additionally, our algorithm achieves a comparable spatial calibration accuracy to Doer et al. [19] on the noisier radar ego-velocity data. However, the high angular velocity trajectory is challenging for real-world camera localization and the amount of excitation is not necessary if the radar data are sufficiently accurate.

## VI. REAL-WORLD EXPERIMENTS

To verify the performance and accuracy of our algorithm, we carried out a series of real-world experiments involving three different radar-camera systems. We discuss the various systems and their implementation details in Section VI-A. In Section VI-B, we show that the set of spatiotemporal calibration parameters estimated by our algorithm have a similar level of alignment accuracy as the parameters estimated by the target-based method of Peršić et al. [2]. In Section VI-C, we demonstrate how spatiotemporal calibration can improve the performance of camera-radar-IMU odometry. Finally, in Section VI-D, we evaluate the accuracy of our algorithm in a challenging situation involving sensors mounted on an autonomous vehicle.

### A. Data Collection and Data Preprocessing

The data collection systems are different for each experiment, but each system includes at least one radar and one camera. The system discussed in Section VI-B is a handheld rig that incorporates a Texas Instruments (TI) AWR1843BOOST radar and Point Grey Flea3 USB camera. The measurement update rates for the sensors are 20 Hz and 30 Hz, respectively. For the experiments in Section VI-C, the data are from the publicly available IRS Radar Thermal Visual Inertial dataset [32]. The data collection system [32] is a handheld rig that mounts on a drone, where measurements are acquired from a TI IWR6843AOP radar, an IDS UI-3241 camera, and an Analog Devices ADIS16448 IMU, operating at frequencies of 10 Hz, 20 Hz, and 409 Hz, respectively. Doer et al. [32] provides additional details about this system. In Section VI-D, the data collection system [33] includes a vehicle-mounted TI AWR1843BOOST radar and three Point Grey Flea3 GigE cameras operating at frequencies of 25 Hz and 16 Hz, respectively.

In our real-world experiments, we use two similar radars that primarily differ in their angular resolutions. If two targets have identical ranges and range-rates but are separated by less than the angular resolution, then the targets will blend together, biasing the radar measurement output. The AWR1843BOOST has azimuth and elevation resolutions of $15°$ and $58°$, respectively, while the IWR6843AOP has azimuth and elevation resolutions of $30°$. As we show in Sections VI-B, VI-C, and VI-D, our algorithm is capable of calibrating both radars even though they have differing angular resolutions. For additional information on the radars used in our experiments, we refer
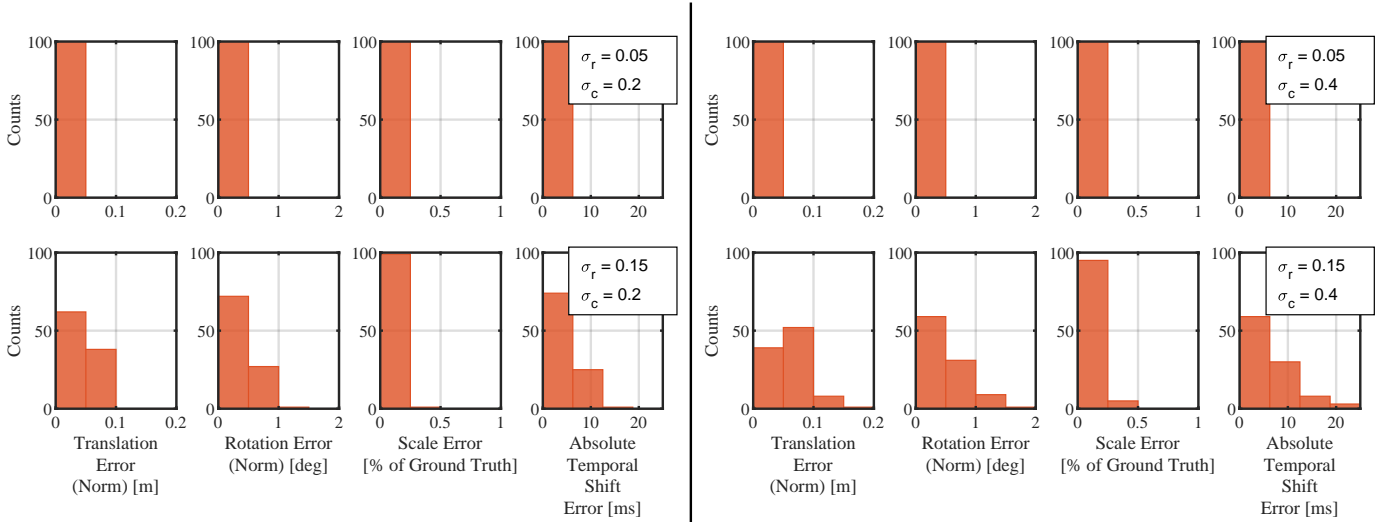
Fig. 5: High linear and low angular velocity trajectory calibration results from our simulation experiments. Each subplot is a histogram of the error between the estimated and true parameter values for 100 trials at a given level of measurement noise. Each row presents the results for a level of measurement noise. The levels of noise are a combination of two radar measurement noise levels ($\sigma_r = 0.05$ or $0.15$ m/s) and two camera pixel measurement noise levels ($\sigma_c = 0.2$ or $0.4$ pixels).
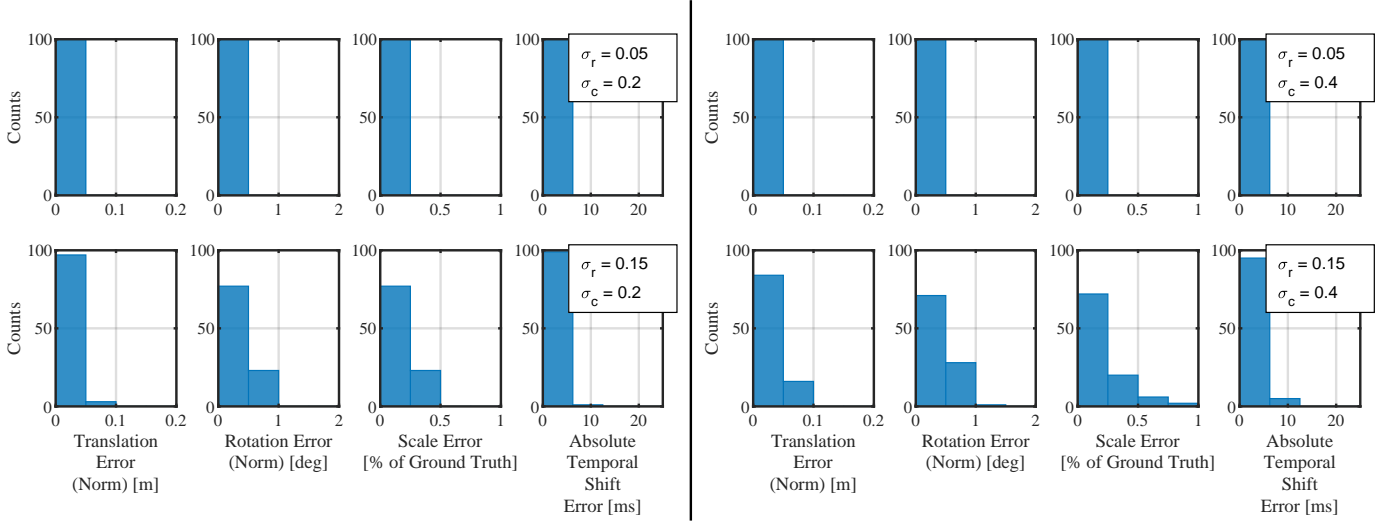


Fig. 6: Low linear and high angular velocity trajectory calibration results from our simulation experiments. Each subplot is a histogram of the error between the estimated and true parameter values for 100 trials at a given level of measurement noise. Each row presents the results for a level of measurement noise. The levels of measurement noise are a combination of two radar measurement noise levels ($\sigma_r = 0.05$ or $0.15$ m/s) and two camera pixel measurement noise levels ($\sigma_c = 0.2$ or $0.4$ pixels).

the reader to the AWR1843BOOST and IWR6843AOP user manuals [34], [35].

To ensure reliable ego-velocity estimation in our experiments, we set the maximum measurable range-rate and constant false alarm rate (CFAR) thresholds for our radar units. The maximum range-rate of the radar must be set above the maximum velocity of the data collection platform because the ego-velocity estimates will saturate at this value. However, an inverse relationship exists between the maximum range-rate and maximum range settings and these must be properly balanced for the operating environment [3]. The on-board radar preprocessing pipeline incorporates a CFAR detector that differentiates targets from background noise in the received EM signal [3]. Since the definition of background noise is also

environment-dependent, we set the CFAR threshold to ensure that the ego-velocity estimator returned a sufficient number of inliers while minimizing the number of outliers. Before each experiment, we performed a series of 'test' data collection runs to tune these settings, making sure that the ego-velocity estimates were not saturating, that there were at least 15 inliers for each measurement, and that the inlier-to-outlier ratio was above 50%.

There are three data preprocessing steps for the experiments discussed in Sections VI-B and VI-C, while the experiment in Section VI-D requires a fourth preprocessing step. Prior to estimating the calibration parameters using our algorithm, we first determine radar ego-velocity estimates using the

algorithm from [25].[2] Second, we rectify the camera images to remove lens distortion effects. Third, we use the feature-based, monocular SLAM algorithm ORB-SLAM3 [36] to provide an initial estimate of the (arbitrarily-scaled) pose of the camera at the time of each image acquisition. While camera pose estimation is possible with any monocular SLAM, we chose this package for its robustness and accuracy [36]. Finally, for the experiment in Section VI-D, we remove outlier radar ego-velocity and camera pose estimates using a median filter. The median filter computes the local median and standard deviation of the signals across a window of time—200 ms and 850 ms for the radar and the camera, respectively. If the measurement at the centre of the window is greater than a chosen threshold from the median, the measurement is treated an outlier. For the tests in Section VI-D, the threshold is set to three standard deviations from the median, since this value eliminates gross outliers without removing noisy, but valid, portions of the signals. We found that this step was necessary to ensure data integrity.

### B. Handheld Rig Experiment

In this experiment, we compared the calibration parameters estimated by our algorithm against the parameters determined by the target-based method in Peršić et al. [2]. To compare the two approaches, we used a handheld rig to collect a dataset consisting of two parts: one part with no visible calibration targets (for our algorithm) and one part with visible targets for target-based calibration. We collected both parts during one continuous run, without power-cycling the sensors. Our quality metric in this case is based on the results from target-based calibration (which can be treated as the 'gold standard,' effectively).

We used the first part of the dataset to perform targetless radar-camera calibration with our algorithm. The procedure consisted of moving the sensor rig, shown in Figure 7, throughout the office environment shown in Figure 9. A segment of the trajectory recovered by our algorithm is plotted in Figure 10. Then, we used the second part of the dataset to perform target-based calibration with the algorithm described in Peršić et al. [2]. In this case, the procedure consisted of moving a trihedral retroreflective target, shown in Figure 8, in front of the stationary radar-camera rig. The second part of the dataset was also used to evaluate the relative accuracy of the parameters estimated by both algorithms.

Our trihedral retroreflective target is specially constructed for calibration evaluation, consisting of a trihedral radar retroreflective 'corner' and a visual AprilTag [37] pattern printed on paper (which is EM-transparent). The target, shown in Figure 8, has the AprilTag mounted in front of the retroreflector. Using the known AprilTag scale, the pose of the camera relative to the AprilTag reference frame can be established. The distance from the origin of the AprilTag frame to the corner of the retroreflector is also known. During data collection, we kept the reflector opening pointed at the radar to ensure a consistent radar reflection.
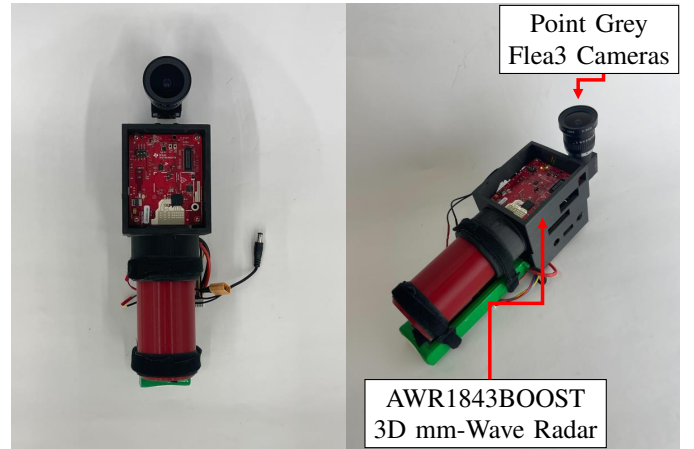


Fig. 7: Two pictures of our handheld sensor rig. The left image is a front view and the right image is an isometric view of the radar-camera unit. The radar antennas are mounted in the white area on the red circuit board. From our CAD model of the handheld rig, the radar-camera translation parameters (i.e., the components of $\mathbf{r}_c^{rc}$) are $r_x = 0.1$, $r_y = 10.5$, and $r_z = -1.0$ cm.
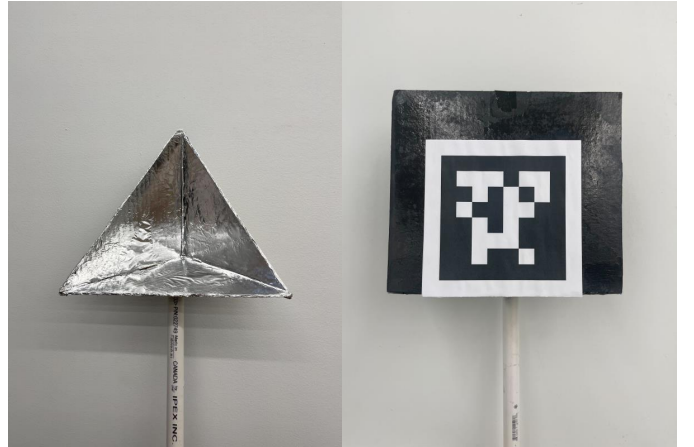


Fig. 8: Our specialized retroreflective radar target used for calibration verification. The left image shows the retroreflector alone, while the right image shows an AprilTag mounted to a flat cardboard backing that is attached to the front of the retroreflector. The cardboard material is fully transparent to the radar EM wave.

We quantify the calibration accuracy based on a 'reprojection error' metric. The reprojection error is the distance between the position of the retroreflector corner predicted from the camera observations and the position measured by the radar, both expressed in the radar frame. The retroreflector is more consistently detected than the AprilTag, and so we linearly interpolate the measured position of the trihedral retroreflector corner to the image timestamps.

Overall, our algorithm achieves results that are comparable to the method from Peršić et al. [2]. Table I shows that the estimated translation and rotation are, per axis, within 1.6 cm and 3 degrees, respectively, of the values estimated by the target-based method. Additionally, our estimated temporal offset differs from the target-based method by only 6 ms. Figure 11 shows that our algorithm, in a completely targetless manner, produces a reprojection error distribution with a

---

[2]Available at: https://github.com/christopherdoer/reve

Fig. 9: Images from our handheld sensor rig calibration dataset, showing two views of the feature-rich indoor test environment.
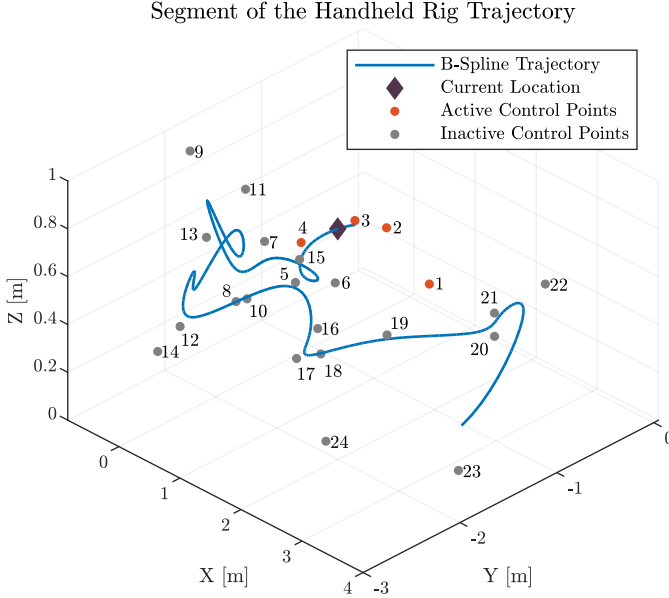


Fig. 10: A segment of the estimated $\mathbf{r}_r^{wr}$ B-spline for the handheld rig during calibration. The purple diamond is position of the rig 6.3 s from the start of the trajectory. The active control points at 6.3 s are shown in orange. As the rig continues along the trajectory, the active control points change.

median that is only 3 mm larger than the target-based method.

### C. IRS Radar Thermal Visual Inertial Datasets

In this section, we demonstrate the versatility of our approach by making use of our estimated calibration parameters to improve the accuracy of camera-radar-IMU odometry. Specifically, we evaluate on the dataset and against the camera-radar-IMU odometry algorithm, known as RRxIO, described by Doer and Trommer in [32]. The extrinsic calibration parameters that accompany the IRS dataset were determined using the radar-IMU extrinsic calibration process described in [19] with ad hoc temporal calibration. Post hoc calibration of the radar and camera is challenging because the test

TABLE I: Calibration parameters for our handheld dataset. The values in each row are estimated by a different algorithm. The rotation between the sensors is given in roll-pitch-yaw (i.e., $\theta_x$, $\theta_y$, $\theta_z$) Euler angle form.

| | $r_x$ [cm] | $r_y$ [cm] | $r_z$ [cm] | $\theta_x$ [rads] | $\theta_y$ [rads] | $\theta_z$ [rads] | $\tau$ [ms] |
|---|---|---|---|---|---|---|---|
| Peršić [2] | -1.60 | 11.9 | -5.02 | -1.59 | 0.07 | -3.12 | -63.8 |
| Ours | -0.48 | 12.2 | -3.42 | -1.62 | 0.02 | -3.15 | -57.9 |

environments do not contain any trihedral reflectors and the motion of the sensor platform is constrained (i.e., there are no deliberate excitations for calibration). To the best of the authors' knowledge, our approach is the only technique that can estimate all of the spatiotemporal calibration parameters for the dataset described in [32].

We chose to calibrate (and to evaluate calibration quality) for three of nine trajectories in the IRS dataset: Gym, MoCap Easy, and MoCap Medium. Data were collected in two environments with varying numbers of features: a large, sparse gymnasium and a feature-rich office setting. For the other six trajectories, poor lighting conditions and rapid motions caused ORB-SLAM3 to fail. To evaluate on a given trajectory, we compute the radar-camera spatiotemporal calibration parameters using our algorithm, and then run RRxIO on the same dataset with our estimated parameters. During evaluation, we disable the live 'camera-to-IMU' extrinsic calibration algorithm that operates as part of RRxIO. Using the known ground truth and the estimated RRxIO trajectories, we are able to determine the quality of our calibration using the following odometric error metrics: the relative translational root mean square error (RMSE RTE), relative rotational RMSE (RRE), absolute translational RMSE (ATE), and absolute rotational RMSE (ARE).

While the parameters estimated by our algorithm, listed in Table II, are relatively close to the parameters provided in Doer and Trommer [32], use of our parameters yields more accurate odometry estimates. The estimated temporal offset for the Gym trajectory is the only large deviation from the values in Doer and Trommer [32]. Table III reports the absolute and relative translation and rotation errors for the RRxIO trajectories after a yaw alignment. The parameters estimated by our algorithm improve the translation error on all datasets and rotation error for two of the datasets. Notably, the Gym dataset, which has the largest temporal offset, improves the most.

### D. Vehicle Experiments

In this section, we verify the accuracy of our calibration algorithm by estimating the distance between cameras mounted on an autonomous vehicle. This task was challenging because,
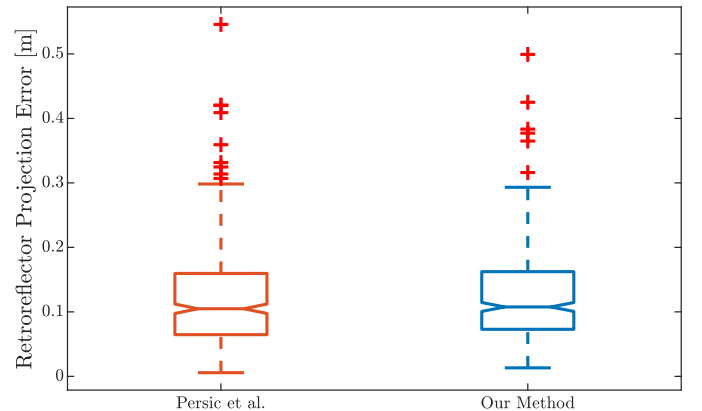


Fig. 11: The reprojection error distributions for the state-of-the-art method in [2] and ours.

TABLE II: Radar-IMU calibration parameters determined for three IRS trajectories and by RRxIO. The radar-IMU calibration parameters listed for our algorithm are a combination of the IRS IMU-camera parameters and our camera-radar parameters. The rotation between the sensors is given by roll-pitch-yaw (i.e., $\theta_x$, $\theta_y$, $\theta_z$) Euler angles.

| | $r_x$ [cm] | $r_y$ [cm] | $r_z$ [cm] | $\theta_x$ [rads] | $\theta_y$ [rads] | $\theta_z$ [rads] | $\tau$ [ms] |
|---|---|---|---|---|---|---|---|
| RRxIO | 6.00 | 4.00 | -4.00 | -3.14 | 0.02 | -1.59 | 8.00 |
| ME† (ours) | 4.08 | 4.71 | -5.05 | -3.12 | 0.01 | -1.59 | 13.1 |
| MM† (ours) | 3.90 | 4.46 | -5.63 | -3.11 | 0.01 | -1.59 | 15.4 |
| Gym (ours) | 3.27 | 4.48 | -3.62 | -3.15 | -0.06 | -1.60 | 40.7 |

† These datasets are MoCap Easy (ME) and MoCap Medium (MM).

as shown in Figure 13, the radar-camera pairs do not share overlapping fields of view, so it is impossible to perform calibration using a target-based method. Additionally, the constrained motion of the car results in a poorly conditioned problem (i.e., the minimum eigenvalue of the identifiability matrix in Equation (24) is close to zero). The poor conditioning of the problem makes the estimated parameters very sensitive to sensor measurement noise, which can lead to inaccurate results. To overcome the poor conditioning of this system, we add an extrinsic calibration prior,

$$\mathbf{e}_{prior} = \log(\mathbf{T}_{cr}^{-1}\mathbf{T}_{cr,prior}),$$

$$J_{prior} = \mathbf{e}_{prior}{}^{T}\boldsymbol{\Sigma}_{prior}^{-1}\mathbf{e}_{prior}, \qquad (25)$$

$$\boldsymbol{\Sigma}_{prior} = \begin{bmatrix} \sigma_t^2\,\mathbf{I}_{3\times3} & \mathbf{0}_{3\times3} \\ \mathbf{0}_{3\times3} & \sigma_\theta^2\,\mathbf{I}_{3\times3} \end{bmatrix},$$

to the optimization problem. For our experiments, the prior for the extrinsic calibration parameters ($\mathbf{T}_{cr,prior}$) is derived from hand measurement. We set the prior uncertainty for the translation ($\sigma_t$) to 0.1 m along each axis, and the prior uncertainty for the rotation to ($\sigma_\theta$) 30 degrees. The addition of this term stabilizes the estimation of the vertical translation between the radar and cameras, in particular. After optimization, less than 1% of the final cost value is due to the prior error term.

The mounting positions of the radar and three cameras on the car are shown in Figure 12, and the corresponding fields of view are shown in Figure 13. The first camera is positioned at the centre of the car and faces the direction of travel. The other two cameras are placed to the left and right of the centre camera and point roughly 45 degrees left and right from the forward axis, respectively. The 3D radar is more than one metre away from all of the cameras, facing towards the rear of the car, opposite the direction of travel.

We collected a total of nine datasets from the radar and the cameras (three datasets per camera) while driving two laps of a figure eight pattern. Data collection took place

TABLE III: Odometry performance evaluation for RRxIO and for our algorithm on three IRS trajectories.

| | RTE [%] | | RRE [deg/m] | | ATE [m] | | ARE [deg] | |
|---|---|---|---|---|---|---|---|---|
| Dataset | RRxIO | Ours | RRxIO | Ours | RRxIO | Ours | RRxIO | Ours |
| ME† | 0.809 | 0.669 | 0.084 | 0.089 | 0.177 | 0.144 | 1.567 | 1.918 |
| MM† | 1.377 | 1.097 | 0.122 | 0.095 | 0.351 | 0.260 | 2.522 | 2.027 |
| Gym | 1.170 | 0.752 | 0.076 | 0.054 | 0.308 | 0.195 | 2.087 | 1.349 |

† These datasets are MoCap easy (ME) and MoCap Medium (MM).



Fig. 12: Two views of the radar and camera mounting positions on the vehicle used in our experiments. The left image shows the mounting position of the TI radar. The right image shows the mounting positions of the three Point Grey cameras. The radar and the cameras do not share overlapping fields of view.

in a sparse parking lot environment, where the radar and camera features were at a substantial distance from the vehicle. We evaluated the accuracy of our estimated parameters by comparing the estimated distances between the centre camera and the two side cameras to the distances measured using a Leica Nova MS50 MultiStation. This method of comparison was selected in part because camera-to-camera extrinsic calibration is difficult for camera pairs that have minimal field of view overlap. Additionally, structural components of the car prevent direct measurement of the distance between the radar and cameras. Each run of our spatiotemporal calibration algorithm produced an estimated extrinsic calibration, for a total of three sets of estimated extrinsic calibration parameters for each camera. The transformations between the centre-to-left and -right cameras are computed by combining two radar extrinsic calibration estimates, which give a total of 18 camera-to-camera extrinsic calibration estimates (nine left and nine right).

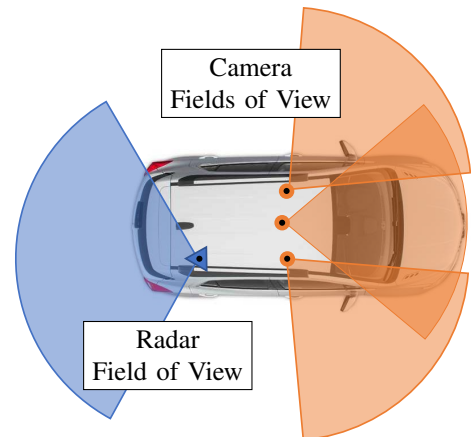Figure 14 shows the distribution of distance errors. The



Fig. 13: Fields of view of the radar and cameras sensors used for our vehicle experiments.

majority of estimated extrinsic calibration parameters result in a camera-to-camera distance error of less than 5 cm, with two values that are greater than 10 cm. This error is reasonable given the 'chained' nature of the two radar-camera transforms. The accuracy of the estimated camera-camera distance depends on the accuracy of the translation and rotation parameters for both transforms. For this experiment, the radar-camera transforms have translation magnitudes greater than 1 m, so a small error in either estimated rotation results in a large distance error. In turn, we expect the estimated rotation and translation parameters to be within $5°$ and 5 cm of the their true values.

### E. Calibration Environment

Several notes are in order regarding environments that are suitable for calibration. Although our algorithm does not require any retroreflective targets for the radar or a specific calibration pattern for the camera, there are nonetheless some limitations on where calibration can be performed. To ensure accurate ego-velocity estimation, the calibration environment should contain, at minimum, four stationary landmarks, with more being better. Also, to ensure accurate camera pose estimation using ORB-SLAM3, the scene should have sufficient lighting and visual texture. As a result, calibration should generally not be performed in scenes with many moving targets, dim lighting, or inclement weather such as fog. The accuracy of the camera pose estimates ultimately depends upon the specific SLAM algorithm that is chosen.

## VII. CONCLUSION

In this paper, we described an algorithm that leverages radar ego-velocity estimates, unscaled camera pose measurements, and a continuous-time trajectory representation to perform targetless radar-to-camera spatiotemporal calibration. We proved that the calibration problem is identifiable and determined the necessary conditions for successful calibration. Through simulation studies, we demonstrated that our algorithm is accurate, but can be sensitive to the amount of noise present in the radar
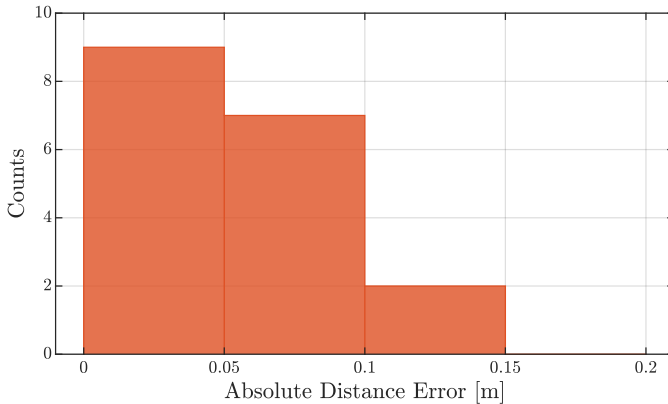


Fig. 14: Results from the vehicle calibration experiment, where the radar and the cameras do not share overlapping fields of view. The distance error is the difference between the estimated and measured distances between the center-left and center-right cameras. The ground truth distance was determined using a Leica MultiStation.

range-rate measurements. Further, we evaluated our algorithm in three different, real-world environments. First, we showed, using data from a handheld sensor rig, that our approach can match the accuracy of target-based calibration methods. Second, we presented results indicating that calibration can improve the localization performance of a hardware-triggered radar-camera-IMU system. Finally, we established that our calibration framework can be applied to AV systems, where the radar and camera are mounted at a significant distance from each other and do not share overlapping fields of view.

There are several potential directions for future research. It would be valuable to develop a method to automatically determine the knot spacing required for the continuous-time spline representation. Our calibration approach could naturally be extended to the multi-camera and multi-radar setting. Other pairs of sensors could also be considered beyond radar-camera pairs, including radar-inertial sensor combinations, for example.

## APPENDIX A
### AN EXTENSION ON THE OBSERVABILITY OF RADAR-TO-CAMERA EXTRINSIC CALIBRATION FROM WISE ET. AL [4]

In this appendix, we provide an extension to our earlier work in [4] demonstrating that radar-to-camera spatial calibration (with a known temporal offset) is locally weakly observable.

### A. Notation for Nonlinear Observability Analysis

In Section III, we represent rotations as orthonormal matrices, which is convenient for the calibration problem but slightly more difficult to use for observability analyses. In this section, we rely on unit quaternions to represent rotations, avoiding the need to use exponential functions. We write a unit quaternion in 'vector' form as

$$\mathbf{q} = \begin{bmatrix} q_0 & \mathbf{q}_v \end{bmatrix}, \tag{26}$$

with the scalar component $q_0$ and vector component $\mathbf{q}_v$, such that $\|\mathbf{q}\|_2 = 1$. The conversion from unit quaternion to rotation matrix is given by the formula

$$\mathbf{R}_{ab} = \mathbf{R}(\mathbf{q}_{ab}) = (2q_0^2 - 1)\mathbf{I}_3 + 2\mathbf{q}_v\,\mathbf{q}_v^T + 2q_0\,\mathbf{q}_v^\wedge. \tag{27}$$

The quaternion kinematics are defined by

$$\dot{\mathbf{q}} = \frac{1}{2}\mathbf{\Xi}(\mathbf{q})\,\boldsymbol{\omega} = \frac{1}{2}\mathbf{\Omega}(\boldsymbol{\omega})\mathbf{q}, \tag{28}$$

$$\mathbf{\Omega}(\boldsymbol{\omega}) = \begin{bmatrix} 0 & -\boldsymbol{\omega}^T \\ \boldsymbol{\omega} & -\boldsymbol{\omega}^\wedge \end{bmatrix}, \tag{29}$$

$$\mathbf{\Xi}(\mathbf{q}) = \begin{bmatrix} -\mathbf{q}_v^T \\ q_0\mathbf{I}_3 + \mathbf{q}_v^\wedge \end{bmatrix}, \tag{30}$$

where $\boldsymbol{\omega}$ is the angular velocity vector. The following identity is useful for the observability proof,

$$\frac{\partial \mathbf{R}(\mathbf{q})\mathbf{p}}{\partial \mathbf{q}} = [(4q_0\mathbf{I}_3 + 2\mathbf{q}_v^\wedge)\mathbf{p} \\ 2((\mathbf{q}_v^T\mathbf{p})\mathbf{I}_3 + \mathbf{q}_v\,\mathbf{p}^T - q_0\mathbf{p}^\wedge)], \tag{31}$$

where $\mathbf{p}$ is an arbitrary $3 \times 1$ vector. If we transpose the rotation matrix on the left side of Equation (31), then the skew-symmetric terms on the right side change sign from positive to negative and vice versa.

### B. Local Weak Observability

Following the procedure outlined in Section IV-A, we define the system equations, compute the respective Lie derivatives, and demonstrate that the nonlinear observability matrix has full column rank. In the analysis here, the pose, velocity, and acceleration states of the radar-camera system are camera-centric (i.e., taken with respect to the camera and not the radar). Since the camera-centric states can be used to determine the radar-centric states, this change does not affect the observability result.

Considering the camera frame $\underset{\rightarrow}{\mathcal{F}}_c$, the radar frame $\underset{\rightarrow}{\mathcal{F}}_r$, and the world frame $\underset{\rightarrow}{\mathcal{F}}_w$, the state vector for the observability analysis is defined as

$$\mathbf{x} = \begin{bmatrix} \mathbf{r}_w^{cwT} & \mathbf{q}_{wc}^{~T} & \mathbf{v}_w^{cwT} & \boldsymbol{\omega}_c^{cwT} & \mathbf{a}_w^{cwT} & \boldsymbol{\alpha}_c^{cwT} \\ \gamma & \mathbf{r}_c^{rcT} & \mathbf{q}_{cr}^{~T} \end{bmatrix}^T, \tag{32}$$

where $\mathbf{r}$, $\mathbf{v}$, and $\mathbf{a}$ denote translation, linear velocity, and linear acceleration, respectively. The vectors $\boldsymbol{\omega}$ and $\boldsymbol{\alpha}$ are the angular velocity and the angular acceleration, respectively. Finally, $\gamma$ is the scale factor for the camera translation (for a monocular camera system). The motion model for the system is

$$\dot{\mathbf{x}} = \mathbf{f}_0(\mathbf{x}) + \mathbf{f}_1(\mathbf{x}) = \begin{bmatrix} \mathbf{0}_{3\times1} \\ \frac{1}{2}\Xi(\mathbf{q}_{wc})\boldsymbol{\omega}_c^{cw} \\ \mathbf{0}_{3\times1} \\ \boldsymbol{\alpha}_c^{cw} \\ \mathbf{0}_{3\times1} \\ \mathbf{0}_{3\times1} \\ 0 \\ \mathbf{0}_{3\times1} \\ \mathbf{0}_{4\times1} \end{bmatrix} + \begin{bmatrix} \mathbf{v}_w^{cw} \\ \mathbf{0}_{4\times1} \\ \mathbf{a}_w^{cw} \\ \mathbf{0}_{3\times1} \\ \mathbf{0}_{3\times1} \\ \mathbf{0}_{3\times1} \\ 0 \\ \mathbf{0}_{3\times1} \\ \mathbf{0}_{4\times1} \end{bmatrix}. \tag{33}$$

The measurement model equations for the (scaled) camera translation and rotation are, respectively,

$$\begin{aligned} \mathbf{h}_1 &= \gamma \, \mathbf{r}_w^{cw}, \\ \mathbf{h}_2 &= \mathbf{q}_{wc}. \end{aligned} \tag{34}$$

Using the camera-centric model, it is possible to directly measure $\mathbf{q}_{wc}$ and, following the result in Section IV-B, to determine $\boldsymbol{\omega}_c^{cw}$ and $\boldsymbol{\alpha}_c^{cw}$. Finally, the radar ego-velocity measurement equation is

$$\mathbf{h}_3 = \mathbf{R}^T(\mathbf{q}_{cr})(\mathbf{R}^T(\mathbf{q}_{wc})\mathbf{v}_w^{cw} + \boldsymbol{\omega}_c^{cw\wedge}\mathbf{r}_c^{rc}). \tag{35}$$

The observability analysis requires the zeroth-, first-, and second-order Lie derivatives. The zeroth-order Lie derivatives are

$$\begin{aligned} \nabla L^0 \mathbf{h}_1 &= \begin{bmatrix} \gamma \mathbf{I}_3 & \mathbf{0}_{3\times16} & \mathbf{r}_w^{cw} & \mathbf{0}_{3\times7} \end{bmatrix}, \\ \nabla L^0 \mathbf{h}_2 &= \begin{bmatrix} \mathbf{0}_{4\times3} & \mathbf{I}_4 & \mathbf{0}_{4\times20} \end{bmatrix}, \\ \nabla L^0 \mathbf{h}_3 &= [\mathbf{0}_{3\times3} \quad \mathbf{A} \quad \mathbf{R}^T(\mathbf{q}_{cr})\mathbf{R}^T(\mathbf{q}_{wc}) \\ &\quad - \mathbf{R}^T(\mathbf{q}_{cr})\mathbf{r}_c^{rc\wedge} \quad \mathbf{0}_{3\times7} \\ &\quad \mathbf{R}^T(\mathbf{q}_{cr})\boldsymbol{\omega}_c^{cw\wedge} \quad \mathbf{B}], \end{aligned} \tag{36}$$

where

$$\begin{aligned} \mathbf{A} &= \mathbf{R}^T(\mathbf{q}_{cr})\frac{\partial \, \mathbf{R}^T(\mathbf{q}_{wc})\mathbf{v}_w^{cw}}{\partial \mathbf{q}_{wc}}, \\ \mathbf{B} &= \frac{\partial \mathbf{R}^T(\mathbf{q}_{cr})(\mathbf{R}^T(\mathbf{q}_{wc})\mathbf{v}_w^{cw} + \boldsymbol{\omega}_c^{cw\wedge}\mathbf{r}_c^{rc})}{\partial \mathbf{q}_{cr}}. \end{aligned} \tag{37}$$

The first-order Lie derivatives are

$$\begin{aligned} \nabla L_{\mathbf{f}_1}^1 \mathbf{h}_1 &= \begin{bmatrix} \mathbf{0}_{3\times7} & \gamma \mathbf{I}_3 & \mathbf{0}_{3\times9} & \mathbf{v}_w^{cw} & \mathbf{0}_{3\times7} \end{bmatrix}, \\ \nabla L_{\mathbf{f}_0}^1 \mathbf{h}_2 &= [\mathbf{0}_{4\times3} \quad \tfrac{1}{2}\Omega(\boldsymbol{\omega}_c^{cw}) \quad \mathbf{0}_{4\times3} \\ &\quad \tfrac{1}{2}\Xi(\mathbf{q}_{wc}) \quad \mathbf{0}_{4\times14}], \\ \nabla L_{\mathbf{f}_0}^1 \mathbf{h}_3 &= [\mathbf{0}_{3\times3} \quad \mathbf{C} \quad \mathbf{D} \quad \mathbf{E} \quad \mathbf{0}_{3\times3} \\ &\quad \mathbf{F} \quad \mathbf{0}_{3\times1} \quad \mathbf{R}^T(\mathbf{q}_{cr})\boldsymbol{\alpha}_c^{cw\wedge} \quad \mathbf{G}], \\ \nabla L_{\mathbf{f}_1}^1 \mathbf{h}_3 &= [\mathbf{0}_{3\times3} \quad \mathbf{H} \quad \mathbf{0}_{3\times6} \\ &\quad \mathbf{R}^T(\mathbf{q}_{cr})\mathbf{R}^T(\mathbf{q}_{wc}) \quad \mathbf{0}_{3\times7} \quad \mathbf{L}], \end{aligned} \tag{38}$$

where

$$\begin{aligned} \mathbf{H} &= \mathbf{R}^T(\mathbf{q}_{cr})\frac{\partial \, \mathbf{R}^T(\mathbf{q}_{wc})\mathbf{a}_w^{cw}}{\partial \mathbf{q}_{wc}}, \\ \mathbf{L} &= \frac{\partial \, \mathbf{R}^T(\mathbf{q}_{cr})\mathbf{R}^T(\mathbf{q}_{wc})\mathbf{a}_w^{cw}}{\partial \mathbf{q}_{cr}}. \end{aligned} \tag{39}$$

We do not explicitly require the nonzero matrices, $\mathbf{C}$, $\mathbf{E}$, and $\mathbf{F}$, in Equation (38) because the submatrix formed from the columns corresponding to the rotation states can be shown to be full rank. The matrices $\mathbf{D}$ and $\mathbf{G}$ are required for the analysis, but we omit them here for brevity. The second-order Lie derivatives are

$$\begin{aligned} \nabla L_{\mathbf{f}_1}^2 \mathbf{h}_1 &= \begin{bmatrix} \mathbf{0}_{3\times13} & \gamma \mathbf{I}_3 & \mathbf{0}_{3\times3} & \mathbf{a}_w^{cw} & \mathbf{0}_{3\times7} \end{bmatrix}, \\ \nabla L_{\mathbf{f}_0}^2 \mathbf{h}_2 &= [\mathbf{0}_{4\times3} \quad \tfrac{1}{4}(2\Omega(\boldsymbol{\alpha}_c^{cw}) - \boldsymbol{\omega}_c^{cwT}\boldsymbol{\omega}_c^{cw}\mathbf{I}_4) \\ &\quad \mathbf{0}_{4\times3} \quad -\tfrac{1}{2}\mathbf{q}_{wc}\boldsymbol{\omega}_c^{cwT} \quad \mathbf{0}_{4\times3} \\ &\quad \tfrac{1}{2}\Xi(\mathbf{q}_{wc}) \quad \mathbf{0}_{4\times8}]. \end{aligned} \tag{40}$$

Stacking the gradients of the Lie derivatives, we arrive at the nonlinear observability matrix,

$$\mathbf{O} = \begin{bmatrix} \nabla L^0 \mathbf{h}_1 \\ \nabla L_{f_1}^1 \mathbf{h}_1 \\ \nabla L_{f_1 f_1}^2 \mathbf{h}_1 \\ \nabla L^0 \mathbf{h}_2 \\ \nabla L_{f_0}^1 \mathbf{h}_2 \\ \nabla L_{f_0 f_0}^2 \mathbf{h}_2 \\ \nabla L^0 \mathbf{h}_3 \\ \nabla L_{f_0}^1 \mathbf{h}_3 \\ \nabla L_{f_1}^1 \mathbf{h}_3 \end{bmatrix}. \tag{41}$$

This matrix can be shown to be full column rank (except when excitation of the system is insufficient), hence the system is locally weakly observable.

## REFERENCES

[1] C.-L. Lee, Y.-H. Hsueh, C.-C. Wang, and W.-C. Lin, "Extrinsic and Temporal Calibration of Automotive Radar and 3D Lidar," in *2020 IEEE/RSJ Intl. Conf. Intelligent Robots and Systems (IROS)*, Las Vegas, NV, USA, Oct. 25, 2020–Jan. 24, 2021 2020, pp. 9976–9983.

[2] J. Peršić, L. Petrović, I. Marković, and I. Petrović, "Spatiotemporal multisensor calibration via Gaussian processes moving target tracking," *IEEE Trans. Robotics*, vol. 37, no. 5, pp. 1401–1415, Mar. 2021.

[3] M. A. Richards, J. A. Scheer, and W. A. Holm, Eds., *Principles of Modern Radar: Basic Principles*. Institution of Eng. and Technol., 2010, vol. 1.

[4] E. Wise, J. Peršić, C. Grebe, I. Petrović, and J. Kelly, "A continuous-time approach for 3D radar-to-camera extrinsic calibration," in *2021 IEEE Intl. Conf. Robotics and Automation (ICRA)*, Xi'an, China, May 30–Jun. 5 2021, pp. 13 164–13 170.

[5] C. C. Stahoviak, "An instantaneous 3D ego-velocity measurement algorithm for frequency modulated continuous wave (FMCW) Doppler radar data," Master's thesis, University of Colorado at Boulder, 2019.

[6] S. Sugimoto, H. Tateda, H. Takahashi, and M. Okutomi, "Obstacle detection using millimeter-wave radar and its visualization on image sequence," in *Int. Conf. Pattern Recognition (ICPR)*, Cambridge, England, Aug. 23–26 2004, pp. 342–345.

[7] T. Wang, N. Zheng, J. Xin, and Z. Ma, "Integrating millimeter wave radar with a monocular vision sensor for on-road obstacle detection applications," *Sensors*, vol. 11, no. 9, pp. 8992–9008, Sep. 2011.

[8] D. Y. Kim and M. Jeon, "Data fusion of radar and image measurements for multi-object tracking via Kalman filtering," *Information Sciences*, vol. 278, pp. 641–652, Sep. 2014.

[9] J. Kim, D. S. Han, and B. Senouci, "Radar and vision sensor fusion for object detection in autonomous vehicle surroundings," in *2018 4th Int. Conf. Ubiquitous and Future Networks (ICUFN)*, Prague, Czech Republic, Jul. 3–6 2018, pp. 76–78.

[10] T. Kim, S. Kim, E. Lee, and M. Park, "Comparative analysis of RADAR-IR sensor fusion methods for object detection," in *2017 17th Int. Conf. Control, Automation and Systems (ICCAS)*, Jeju, Korea, Oct. 18–21 2017, pp. 1576–1580.

[11] G. El Natour, O. Ait Aider, R. Rouveure, F. Berry, and P. Faure, "Radar and vision sensors calibration for outdoor 3D reconstruction," in *2015 IEEE Int. Conf. Robotics and Automation (ICRA)*, Seattle, WA, USA, May 25–30 2015, pp. 2084–2089.

[12] J. Domhof, J. F. P. Kooij, and D. M. Gavrila, "An extrinsic calibration tool for radar, camera and lidar," in *2019 Int. Conf. Robotics and Automation (ICRA)*, Montréal, Canada, May, 20–24 2019, pp. 8107–8113.

[13] J. Peršić, I. Marković, and I. Petrović, "Extrinsic 6DoF calibration of a radar–lidar–camera system enhanced by radar cross section estimates evaluation," *Robotics and Autonomous Systems*, vol. 114, pp. 217–230, Apr. 2019.

[14] J. Oh, K. Kim, M. Park, and S. Kim, "A comparative study on camera-radar calibration methods," in *2018 15th Int. Conf. Control, Automation, Robotics and Vision (ICARCV)*, Singapore, Nov. 18–21 2018, pp. 1057–1062.

[15] C. Schöller, M. Schnettler, A. Krämmer, G. Hinz, M. Bakovic, M. Güzet, and A. Knoll, "Targetless rotational auto-calibration of radar and camera for intelligent transportation systems," in *2019 IEEE Intelligent Transportation Systems Conf. (ITSC)*, Auckland, New Zealand, Oct. 27–30 2019, pp. 3934–3941.

[16] J. Peršić, L. Petrović, I. Marković, and I. Petrović, "Online multi-sensor calibration based on moving object tracking," *Advanced Robotics*, vol. 35, no. 3–4, pp. 130–140, Sep. 2021.

[17] L. Heng, "Automatic targetless extrinsic calibration of multiple 3D lidars and radars," in *2020 IEEE/RSJ Intl. Conf. Intelligent Robots and Systems (IROS)*, Las Vegas, NV, USA, Oct. 25, 2020–Jan. 24, 2021 2020, pp. 10 669–10 675.

[18] D. Kellner, M. Barjenbruch, K. Dietmayer, J. Klappstein, and J. Dickmann, "Joint radar alignment and odometry calibration," in *2015 18th Int. Conf. Information Fusion (FUSION)*, Washington, DC, USA, Jul. 6–9 2015, pp. 366–374.

[19] C. Doer and G. F. Trommer, "Radar inertial odometry with online calibration," in *2020 European Navigation Conf. (ENC)*, Nov. 23–24 2020, pp. 1–10.

[20] J. Rehder, R. Siegwart, and P. Furgale, "A general approach to spatiotemporal calibration in multisensor systems," *IEEE Trans. Robotics*, vol. 32, no. 2, pp. 383–398, Apr. 2016.

[21] T. D. Barfoot, *State estimation for robotics*. Cambridge, UK: Cambridge Univ. Press, 2017.

[22] C. Sommer, V. Usenko, D. Schubert, N. Demmel, and D. Cremers, "Efficient derivative computation for cumulative B-splines on Lie groups," in *2020 IEEE/CVF Conf. Computer Vision and Pattern Recognition (CVPR)*, Jun. 14–19 2020, pp. 11 145–11 153.

[23] C. de Boor, *A Practical Guide to Splines*, ser. Applied Mathematical Sciences. Springer-Verlag, Jan. 1978, vol. 27.

[24] K. Qin, "General matrix representations for b-splines," in *6th Pacific Conf. on Computer Graphics and Applications*, Singapore, Oct. 26–29 1998, pp. 37–43.

[25] C. Doer and G. F. Trommer, "An EKF based approach to radar inertial odometry," in *2020 IEEE Intl. Conf. Multisensor Fusion and Integration for Intelligent Systems (MFI)*, Karlsruhe, Germany, Sep. 14–16 2020, pp. 152–159.

[26] A. Chiuso, P. Favaro, Hailin Jin, and S. Soatto, "Structure from motion causally integrated over time," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 24, no. 4, pp. 523–535, Apr. 2002.

[27] S. Agarwal *et al.*, *Ceres Solver*. [Online]. Available: http://ceres-solver.org

[28] M. Li and A. I. Mourikis, "Online temporal calibration for camera-IMU systems: Theory and algorithms," *Intl. J. Robotics Research*, vol. 33, no. 7, pp. 947–964, 2014.

[29] R. A. Hewitt and J. A. Marshall, "Towards intensity-augmented SLAM with LiDAR and ToF sensors," in *Proc. IEEE/RSJ Intl. Conf. Intelligent Robots and Systems (IROS)*, Hamburg, Germany, September/October 2015, pp. 1957–1961.

[30] R. Hermann and A. Krener, "Nonlinear controllability and observability," *IEEE Trans. Automatic Control*, vol. 22, no. 5, pp. 728–740, Oct. 1977.

[31] J. Kelly, C. Grebe, and M. Giamou, "A question of time: Revisiting the use of recursive filtering for temporal calibration of multisensor systems," in *Proc. IEEE Intl. Conf. Multisensor Fusion and Integration (MFI)*, Karlsruhe, Germany, 2021.

[32] C. Doer and G. F. Trommer, "Radar visual inertial odometry and radar thermal inertial odometry: Robust navigation even in challenging visual conditions," in *2021 IEEE/RSJ Intl. Conf. Intelligent Robots and Systems (IROS)*, Prague, Czech Republic, Sep. 27 – Oct. 1 2021, pp. 331–338.

[33] K. Burnett, D. J. Yoon, Y. Wu, A. Z. Li, H. Zhang, S. Lu, J. Qian, W.-K. Tseng, A. Lambert, K. Y. Leung, A. P. Schoellig, and T. D. Barfoot, "Boreas: A multi-season autonomous driving dataset," *The International Journal of Robotics Research*, vol. 42, pp. 33–42, 2023.

[34] Texas Instruments, *xWR1843 Evaluation Module (xWR1843BOOST) Single-Chip mmWave Sensing Solution*, May 2020. [Online]. Available: https://www.ti.com/lit/ug/spruim4b/spruim4b.pdf

[35] ——, *IWR6843AOP Single-Chip 60- to 64-GHz mmWave Sensor Antennas-On-Package (AOP)*, July 2022. [Online]. Available: https://www.ti.com/lit/ds/symlink/iwr6843aop.pdf

[36] C. Campos, R. Elvira, J. J. G. Rodr'iguez, J. M. M. Montiel, and J. D. Tardós, "ORB-SLAM3: An accurate open-source library for visual, visual–inertial, and multimap SLAM," *IEEE Trans. Robotics*, vol. 37, no. 6, pp. 1874–1890, May 2021.

[37] E. Olson, "AprilTag: A robust and flexible visual fiducial system," in *2011 IEEE Intl. Conf. Robotics and Automation (ICRA)*, Shanghai, China, May 9–13 2011, pp. 3400–3407.

**Emmett Wise** received his Bachelor of Applied Science in engineering physics from Queens University, Kingston, Canada, in 2016. After graduation, he worked to automate the manufacturing of nanoscale coatings at 3M Canada. He is currently a Ph.D. Candidate in the Space and Terrestrial Autonomous Robotics (STARS) Laboratory at the Institute for Aerospace Studies, Toronto, Canada. His research interests include perception, calibration, and state estimation.

**Qilong (Jerry) Cheng** is currently pursuing his Master of Engineering degree in electrical and computer engineering, having completed a Bachelor's in Mechanical Engineering at the University of Toronto. He was a developer at Autodesk from 2019 to 2020. From 2020 to 2021, he made design patent contributions to a high pressure spray nozzle design for the China State Shipbuilding Corporation. He is currently a graduate research student in the Space and Terrestrial Autonomous Robotics (STARS) Laboratory at the University of Toronto, focusing on calibration and state estimation.

**Jonathan Kelly** received the Ph.D. degree in Computer Science from the University of Southern California, Los Angeles, USA, in 2011. From 2011 to 2013 he was a postdoctoral associate in the Computer Science and Artificial Intelligence Laboratory at the Massachusetts Institute of Technology, Cambridge, USA. He is currently an associate professor and director of the Space and Terrestrial Autonomous Robotic Systems (STARS) Laboratory at the University of Toronto Institute for Aerospace Studies, Toronto, Canada. Prof. Kelly holds the Tier II Canada Research Chair in Collaborative Robotics. His research interests include perception, planning, and learning for interactive robotic systems.