

Inverse Depth Parametrization for Monocular SLAM

Javier Civera, Andrew J. Davison, J.M.M Montiel, *Member, IEEE*

Abstract— We present a new parametrization for point features within monocular SLAM which permits efficient and accurate representation of uncertainty during undelayed initialisation and beyond, all within the standard EKF (Extended Kalman Filter). The key concept is direct parametrization of the inverse depth of features relative to the camera locations from which they were first viewed, which produces measurement equations with a high degree of linearity. Importantly, our parametrization can cope with features over a huge range of depths, even those which are so far from the camera that they present little parallax during motion — maintaining sufficient representative uncertainty that these points retain the opportunity to ‘come in’ smoothly from infinity if the camera makes larger movements. Feature initialization is undelayed in the sense that even distant features are immediately used to improve camera motion estimates, acting initially as bearing references but not permanently labelled as such.

The inverse depth parametrization remains well behaved for features at all stages of SLAM processing, but has the drawback in computational terms that each point is represented by a six dimensional state vector as opposed to the standard three of a Euclidean XYZ representation. We show that once the depth estimate of a feature is sufficiently accurate, its representation can safely be converted to the Euclidean XYZ form, and propose a linearity index which allows automatic detection and conversion to maintain maximum efficiency — only low parallax features need be maintained in inverse depth form for long periods.

We present a real-time implementation at 30Hz where the parametrization is validated in a fully automatic 3D SLAM system featuring a hand-held single camera with no additional sensing. Experiments show robust operation in challenging indoor and outdoor environments with very large ranges of scene depth, varied motion and also real-time 360° loop closing.

Index Terms— Real-time vision, monocular SLAM.

I. INTRODUCTION

A monocular camera is a projective sensor which measures the bearing of image features. Given an image sequence of a rigid 3D scene taken from a moving camera, it is now well known that it is possible to compute both the scene structure and the camera motion up to a scale factor. To infer the 3D position of each feature, the moving camera must observe it repeatedly, each time capturing a ray of light from the feature to its optic center. The measured angle between the captured

rays from different viewpoints is the feature’s *parallax* — this is what allows its depth to be estimated.

In off-line ‘Structure from Motion (SfM)’ solutions from the computer vision literature (e.g. [11], [23]), motion and structure are estimated from an image sequence by first applying robust feature matching between pairs or other short overlapping sets of images to estimate relative motion. An optimization procedure then iteratively refines global camera location and scene feature position estimates such that features project as closely as possible to their measured image positions (bundle adjustment). Recently work in the spirit of these methods but with ‘sliding window’ processing and refinement rather than global optimization has produced impressive real-time ‘Visual Odometry’ results when applied to stereo sequences in [21] and for monocular sequences in [20].

An alternative approach to achieving real-time motion and structure estimation are on-line visual SLAM (Simultaneous Localization And Mapping) approaches which use a probabilistic filtering approach to sequentially update estimates of the positions of features (the map) and the current location of the camera. These SLAM methods have different strengths and weaknesses to visual odometry, being able to build consistent and drift-free global maps but with a bounded number of mapped features. The core single Extended Kalman Filter (EKF) SLAM technique, previously proven in multi-sensor robotic applications, was first applied successfully to real-time monocular camera tracking by Davison et al. [8] [9] in a system which built sparse room-sized maps at 30Hz.

A significant limitation of Davison’s and similar approaches, however, was that they could only make use of features which were close to the camera relative to its distance of translation, and therefore exhibited significant parallax during motion. The problem was in initialising uncertain depth estimates for distant features: in the straightforward Euclidean XYZ feature parametrization adopted, position uncertainties for low parallax features are not well represented by the Gaussian distributions implicit in the EKF. The depth coordinate of such features has a probability density which rises sharply at a well-defined minimum depth to a peak, but then tails off very slowly towards infinity — from low parallax measurements it is very difficult to tell whether a feature has a depth of 10 units rather than 100, 1000 or more. For the rest of the paper we refer to Euclidean XYZ parametrization simply as XYZ.

There have been several recent methods proposed for coping with this problem, relying on generally undesirable special treatment of newly initialized features. In this paper we describe a new feature parametrization which is able smoothly to cope with initialization of features at all depths — even up to ‘infinity’ — within the standard EKF framework. The key

Manuscript received February 27, 2007; revised December 6, 2007. This work was supported by Spanish PR2007-0427, DPI2006-13578, DGA(CONSI+D)-CAI IT12-06, EPSRC grant GR/T24685, an EPSRC Advanced Research Fellowship to AJD, Royal Society International Joint Project grant between the University of Oxford, University of Zaragoza and Imperial College London and RAWSEEDS FP6-IST-045144.

J. Civera and J.M.M Montiel are with University of Zaragoza, Dpto. Informtica. Mara de Luna 1, 50018 Zaragoza (Spain)(e-mail:josemari@unizar.es; jcivera@unizar.es)

A.J. Davison is with Imperial College London, Department of Computing, 180 Queen’s Gate SW7 2AZ, UK. (e-mail:ajd@doc.ic.ac.uk)

concept is direct parametrization of inverse depth relative to the camera position from which a feature was first observed.

A. Delayed and Undelayed Initialization

The most obvious approach to coping with feature initialization within a monocular SLAM system is to treat newly detected features separately from the main map, accumulating information in special processing over several frames to reduce depth uncertainty before insertion into the full filter with a standard XYZ representation. Such *delayed initialization* schemes (e.g. [8], [14], [3]) have the drawback that new features, held outside the main probabilistic state, are not able to contribute to the estimation of the camera position until finally included in the map. Further, features which retain low parallax over many frames (those very far from the camera, or close to the motion epipole) are usually rejected completely because they never pass the test for inclusion.

In the delayed approach of Bailey [2], initialization is delayed until the measurement equation is approximately Gaussian and the point can be safely triangulated; here the problem was posed in 2D and validated in simulation. A similar approach for 3D monocular vision with inertial sensing was proposed in [3]. Davison [8] reacted to the detection of a new feature by inserting a 3D semi-infinite ray into the main map representing everything about the feature except its depth, and then used an auxiliary particle filter to explicitly refine the depth estimate over several frames, taking advantage of all the measurements in a high frame-rate sequence but again with new features held outside the main state vector until inclusion.

More recently, several *undelayed initialization* schemes have been proposed, which still treat new features in a special way but are able to benefit immediately from them to improve camera motion estimates — the key insight being that while features with highly uncertain depths provide little information on camera translation, they are extremely useful as bearing references for orientation estimation. The undelayed method proposed by Kwok and Dissanayake [15] was a multiple hypothesis scheme, initializing features at various depths and pruning those not reobserved in subsequent images.

Sola et al. [25], [24] described a more rigorous undelayed approach using a Gaussian Sum Filter approximated by a Federated Information Sharing method to keep the computational overhead low. An important insight was to spread the Gaussian depth hypotheses along the ray according to inverse depth, achieving much better representational efficiency in this way. This method can perhaps be seen as the direct stepping stone between Davison’s particle method and our new inverse depth scheme; a Gaussian sum is a more efficient representation than particles (efficient enough that the separate Gaussians can all be put into the main state vector), but not as efficient as the single Gaussian representation that the inverse depth parametrization allows. Note that neither [15] nor [25] consider features at very large ‘infinite’ depths.

B. Points at Infinity

A major motivation of the approach in this paper is not only efficient undelayed initialization, but also the desire to cope

with features at *all* depths, particularly in outdoor scenes. In SFM, the well-known concept of a point at infinity is a feature which exhibits no parallax during camera motion due to its extreme depth. A star for instance would be observed at the same image location by a camera which translated through many kilometers pointed up at the sky without rotating. Such a feature cannot be used for estimating camera translation but is a perfect bearing reference for estimating rotation. The homogeneous coordinate systems of visual projective geometry used normally in SFM allow explicit representation of points at infinity, and they have proven to play an important role during off-line structure and motion estimation.

In a sequential SLAM system, the difficulty is that we do not know in advance which features are infinite and which are not. Montiel and Davison [19] showed that in the special case where *all features are known to be infinite* — in very large scale outdoor scenes or when the camera rotates on a tripod — SLAM in pure angular coordinates turns the camera into a real-time visual compass. In the more general case, let us imagine a camera moving through a 3D scene with observable features at a range of depths. From the estimation point of view, we can think of all features starting at infinity and ‘coming in’ as the camera moves far enough to measure sufficient parallax. For nearby indoor features, only a few centimetres of movement will be sufficient. Distant features may require many meters or even kilometers of motion before parallax is observed. It is important that these features are not permanently labelled as infinite — a feature that seems to be at infinity should always have the chance to prove its finite depth given enough motion, or there will be the serious risk of systematic errors in the scene map. Our probabilistic SLAM algorithm must be able to represent the uncertainty in depth of seemingly infinite features. Observing no parallax for a feature after 10 units of camera translation does tell us something about its depth — it gives a reliable lower bound, which depends on the amount of motion made by the camera (if the feature had been closer than this we *would* have observed parallax). This explicit consideration of uncertainty in the locations of points has not been previously required in off-line computer vision algorithms, but is very important in the more difficult on-line case.

C. Inverse Depth Representation

Our contribution is to show that in fact there is a unified and straightforward parametrization for feature locations which can handle both initialisation and standard tracking of both close and very distant features within the standard EKF framework. An explicit parametrization of the *inverse depth* of a feature along a semi-infinite ray from the position from which it was first viewed allows a Gaussian distribution to cover uncertainty in depth which spans a depth range from nearby to infinity, and permits seamless crossing over to finite depth estimates of features which have been apparently infinite for long periods of time. The unified representation means that our algorithm requires no special initialisation process for features. They are simply tracked right from the start, immediately contribute to improved camera estimates and have their correlations with

all other features in the map correctly modelled. Note that our parameterization would be equally compatible with other variants of Gaussian filtering such as sparse information filters.

We introduce a linearity index and use it to analyze and prove the representational capability of the inverse depth parametrization for both low and high-parallax features. The only drawback of the inverse depth scheme is the computational issue of increased state vector size, since an inverse depth point needs six parameters rather than the three of XYZ coding. As a solution to this, we show that our linearity index can also be applied to the XYZ parametrization to signal when a feature can be safely switched from inverse depth to XYZ; the usage of the inverse depth representation can in this way be restricted to low parallax feature cases where the XYZ encoding departs from Gaussianity. Note that this ‘switching’, unlike in delayed initialization methods, is purely to reduce computational load; SLAM accuracy with or without switching is almost the same.

The fact is that the projective nature of a camera means that the image measurement process is nearly linear in this inverse depth coordinate. Inverse depth is a concept used widely in computer vision: it appears in the relation between image disparity and point depth in stereo vision; it is interpreted as the parallax with respect to the plane at infinity in [12]. Inverse depth is also used to relate the motion field induced by scene points with the camera velocity in optical flow analysis [13]. In the tracking community, ‘modified polar coordinates’ [1] also exploit the linearity properties of the inverse depth representation in the slightly different, but closely related, problem of target motion analysis (TMA) from measurements gathered by a bearing-only sensor with known motion.

However, the inverse depth idea has not previously been properly integrated in sequential, probabilistic estimation of motion and structure. It has been used in EKF based sequential depth estimation from camera known motion [16] and in multi-baseline stereo Okutomi and Kanade [22] used the inverse depth to increase matching robustness for scene symmetries; matching scores coming from multiple stereo pairs with different baselines were accumulated in a common reference coded in inverse depth, this paper focusing on matching robustness and not on probabilistic uncertainty propagation. In [5] Chowdhury and Chellappa proposed a sequential EKF process using inverse depth but this was some way short of full SLAM in its details. Images are first processed pairwise to obtain a sequence of 3D motions which are then fused with an individual EKF per feature.

It is our parametrization of inverse depth *relative to the positions from which features were first observed* which means that a Gaussian representation is uniquely well behaved, and this is the reason why a straightforward parametrization of monocular SLAM in the homogeneous coordinates of SFM will not give a good result — that representation only meaningfully represents points which appear to be infinite relative to the coordinate origin. It could be said in projective terms that our method defines separate but correlated projective frames for each feature. Another interesting comparison is between our method, where the representation for each feature includes the camera position from which it was first observed and

smoothing/Full SLAM schemes where all historical sensor pose estimates are maintained in a filter.

Two recently published papers from other authors have developed methods which are quite similar to ours. Trawny and Roumeliotis in [26] proposed an undelayed initialization for 2D monocular SLAM which encodes a map point as the intersection of two projection rays. This representation is overparametrized but allows undelayed initialization and encoding of both close and distant features, the approach validated with simulation results.

Eade and Drummond presented an inverse depth initialization scheme within the context of their FastSLAM-based system for monocular SLAM [10], offering some of the same arguments about advantages in linearity as in our paper. The position of each new partially initialised feature added to the map is parametrized with three coordinates representing its direction and inverse depth relative to the camera pose at the first observation, and estimates of these coordinates are refined within a set of Kalman Filters for each particle of the map. Once the inverse depth estimation has collapsed, the feature is converted to a fully initialised standard XYZ representation. While retaining the differentiation between partially and fully-initialised features, they go further and are able to use measurements of partially initialised features with unknown depth to improve estimates of camera orientation and translation via a special epipolar update step. Their approach certainly appears appropriate within a FastSLAM implementation. However, it lacks the satisfying unified quality of the parametrization we present in this paper, where the transition from partially to fully initialised need not be explicitly tackled and full use is automatically made of all of the information available in measurements.

This paper offers a comprehensive and extended version of our work previously published as two conference papers [18] [7]. We now present a full real-time implementation of the inverse depth parameterization which can map up to 50-70 features in real-time on a standard laptop computer. Experimental validation has shown the important role of an accurate camera calibration to improve the system performance especially with wide angle cameras. Our results section includes new real-time experiments, including the key result of vision-only loop closing. Input test image sequences and movies showing the computed solution are included in the paper as multimedia material.

Section II is devoted to defining the state vector, including the camera motion model, XYZ point coding and inverse depth point parametrization. The measurement equation is described in Section III. Section IV presents a discussion about measurement equation linearization errors. Next, feature initialization from a single feature observation is detailed in Section V. In Section VI the switch from inverse depth to XYZ coding is presented, and in Section VII we present experimental validations over real image sequences captured at 30Hz in large scale environments, indoors and outdoors, including real-time performance and a loop closing experiment; links to movies showing the system performance are provided. Finally Section VIII is devoted to conclusions.

II. STATE VECTOR DEFINITION

A. Camera Motion

A constant angular and linear velocity model is used to model hand-held camera motion. The camera state \mathbf{x}_v is composed of pose terms: \mathbf{r}^{WC} camera optical center position and \mathbf{q}^{WC} quaternion defining orientation; and linear and angular velocity \mathbf{v}^W and ω^C relative to world frame W and camera frame C .

We assume that linear and angular accelerations \mathbf{a}^W and α^C affect the camera, producing at each step an impulse of linear velocity, $\mathbf{V}^W = \mathbf{a}^W \Delta t$, and angular velocity $\Omega^C = \alpha^C \Delta t$, with zero mean and known Gaussian distribution. We currently assume a diagonal covariance matrix for the unknown input linear and angular accelerations.

The state update equation for the camera is:

$$\mathbf{f}_v = \begin{pmatrix} \mathbf{r}_{k+1}^{WC} \\ \mathbf{q}_{k+1}^{WC} \\ \mathbf{v}_{k+1}^W \\ \omega_{k+1}^C \end{pmatrix} = \begin{pmatrix} \mathbf{r}_k^{WC} + (\mathbf{v}_k^W + \mathbf{V}_k^W) \Delta t \\ \mathbf{q}_k^{WC} \times \mathbf{q}((\omega_k^C + \Omega^C) \Delta t) \\ \mathbf{v}_k^W + \mathbf{V}_k^W \\ \omega_k^C + \Omega^C \end{pmatrix}, \quad (1)$$

where $\mathbf{q}((\omega_k^C + \Omega^C) \Delta t)$ is the quaternion defined by the rotation vector $(\omega_k^C + \Omega^C) \Delta t$.

B. Euclidean XYZ Point Parametrization

The standard representation for scene points i in terms of Euclidean XYZ coordinates (see Fig 1) is:

$$\mathbf{x}_i = (X_i \ Y_i \ Z_i)^\top. \quad (2)$$

In the paper we refer to the Euclidean XYZ coding simply as XYZ coding.

C. Inverse Depth Point Parametrization

In our new scheme, a scene 3D point i can be defined by the dimension 6 state vector:

$$\mathbf{y}_i = (x_i \ y_i \ z_i \ \theta_i \ \phi_i \ \rho_i)^\top, \quad (3)$$

which models a 3D point located at (see Fig 1):

$$\mathbf{x}_i = \begin{pmatrix} X_i \\ Y_i \\ Z_i \end{pmatrix} = \begin{pmatrix} x_i \\ y_i \\ z_i \end{pmatrix} + \frac{1}{\rho_i} \mathbf{m}(\theta_i, \phi_i) \quad (4)$$

$$\mathbf{m} = (\cos \phi_i \sin \theta_i, -\sin \phi_i, \cos \phi_i \cos \theta_i)^\top. \quad (5)$$

The \mathbf{y}_i vector encodes the ray from the first camera position from which the feature was observed by x_i, y_i, z_i , the camera optical center, and θ_i, ϕ_i azimuth and elevation (coded in the world frame) defining unit directional vector $\mathbf{m}(\theta_i, \phi_i)$. The point's depth along the ray d_i is encoded by its inverse $\rho_i = 1/d_i$.

D. Full State Vector

As in standard EKF SLAM, we use a single joint state vector containing camera pose and feature estimates, with the assumption that the camera moves with respect to a static scene. The whole state vector \mathbf{x} is composed of the camera and all the map features:

$$\mathbf{x} = (\mathbf{x}_v^\top, \mathbf{y}_1^\top, \mathbf{y}_2^\top, \dots, \mathbf{y}_n^\top)^\top. \quad (6)$$

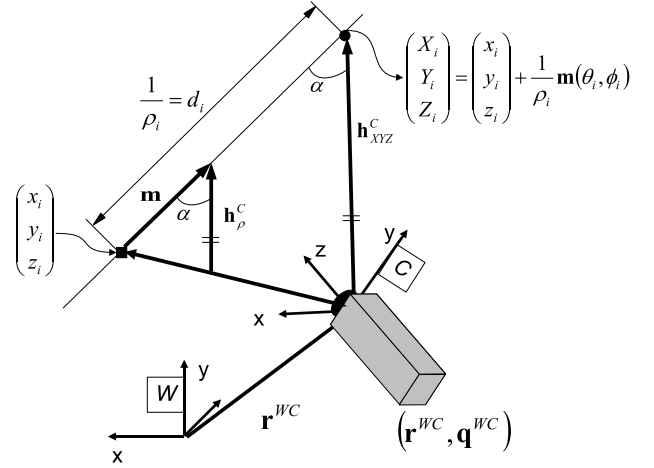


Fig. 1. Feature parametrization and measurement equation.

III. MEASUREMENT EQUATION

Each observed feature imposes a constraint between the camera location and the corresponding map feature (see Fig 1). Observation of a point \mathbf{y}_i (x_i, y_i, z_i) defines a ray coded by a directional vector in the camera frame $\mathbf{h}^C = (h_x \ h_y \ h_z)^\top$. For points in XYZ:

$$\mathbf{h}^C = \mathbf{h}_{XYZ}^C = \mathbf{R}^{CW} \begin{pmatrix} X_i \\ Y_i \\ Z_i \end{pmatrix} - \mathbf{r}^{WC}. \quad (7)$$

For points in inverse depth:

$$\mathbf{h}^C = \mathbf{h}_\rho^C = \mathbf{R}^{CW} \left(\rho_i \begin{pmatrix} x_i \\ y_i \\ z_i \end{pmatrix} - \mathbf{r}^{WC} \right) + \mathbf{m}(\theta_i, \phi_i), \quad (8)$$

where the directional vector has been normalized using the inverse depth. It is worth noting that (8) can be safely used even for points at infinity i.e $\rho_i = 0$.

The camera does not directly observe \mathbf{h}^C but its projection in the image according to the pinhole model. Projection to a normalized retina and then camera calibration is applied:

$$\mathbf{h} = \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} u_0 - \frac{f}{d_x} \frac{h_x}{h_z} \\ v_0 - \frac{f}{d_y} \frac{h_y}{h_z} \end{pmatrix}, \quad (9)$$

where u_0, v_0 is the camera's principal point, f is the focal length and d_x, d_y the pixel size. Finally, a distortion model has to be applied to deal with real camera lenses. In this work we have used the standard two parameter distortion model from photogrammetry [17] (see Appendix for details).

It is worth noting that the measurement equation in inverse depth has a sensitive dependency on the parallax angle α (see Figure 1). At low parallax, Equation (8) can be approximated by $\mathbf{h}^C \approx \mathbf{R}^{CW} (\mathbf{m}(\theta_i, \phi_i))$, and hence the measurement equation only provides information about the camera orientation and the directional vector $\mathbf{m}(\theta_i, \phi_i)$.

IV. MEASUREMENT EQUATION LINEARITY

The more linear the measurement equation is, the better a Kalman Filter performs. This section is devoted to presenting

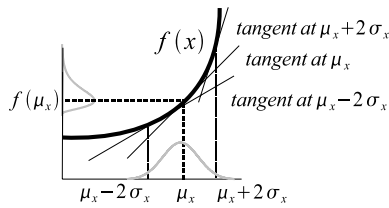


Fig. 2. The first derivative variation in $[\mu_x - 2\sigma_x, \mu_x + 2\sigma_x]$ codes the departure from Gaussianity in the propagation of the uncertain variable through a function.

an analysis of measurement equation linearity for both XYZ and inverse depth codings. These linearity analyses theoretically support the superiority of the inverse depth coding.

A. Linearized propagation of a Gaussian

Let x be an uncertain variable with Gaussian distribution $x \sim N(\mu_x, \sigma_x^2)$. The transformation of x through the function f is a variable y which can be approximated with Gaussian distribution:

$$y \sim N(\mu_y, \sigma_y^2), \quad \mu_y = f(\mu_x), \quad \sigma_y^2 = \left. \frac{\partial f}{\partial x} \right|_{\mu_x} \sigma_x^2 \left. \frac{\partial f}{\partial x} \right|_{\mu_x}^T, \quad (10)$$

if the function f is linear in an interval around μ_x (Figure 2). The interval size in which the function has to be linear depends on σ_x ; the bigger σ_x the wider the interval has to be to cover a significant fraction of the random variable x values. In this work we fix the linearity interval to the 95% confidence region defined by $[\mu_x - 2\sigma_x, \mu_x + 2\sigma_x]$.

If a function is linear in an interval, the first derivative is constant in that interval. To analyze the first derivative variation around the interval $[\mu_x - 2\sigma_x, \mu_x + 2\sigma_x]$ consider the Taylor expansion for the *first derivative*:

$$\left. \frac{\partial f}{\partial x} \right|_{\mu_x + \Delta x} \approx \left. \frac{\partial f}{\partial x} \right|_{\mu_x} + \left. \frac{\partial^2 f}{\partial x^2} \right|_{\mu_x} \Delta x. \quad (11)$$

We propose to compare the value of the derivative at the interval center, μ_x , with the value at the extremes $\mu_x \pm 2\sigma_x$, where the deviation from linearity will be maximal, using the following dimensionless linearity index:

$$L = \left| \frac{\left. \frac{\partial^2 f}{\partial x^2} \right|_{\mu_x} 2\sigma_x}{\left. \frac{\partial f}{\partial x} \right|_{\mu_x}} \right|. \quad (12)$$

When $L \approx 0$, the function can be considered linear in the interval, and hence Gaussianity is preserved during transformation.

B. Linearity of XYZ Parametrization

The linearity of the XYZ representation is analyzed by means of a simplified model which only estimates the depth of a point with respect to the camera. In our analysis, a scene point is observed by two cameras (Figure 3a), both of which are oriented towards the point. The first camera detects the ray on which the point lies. The second camera observes the same

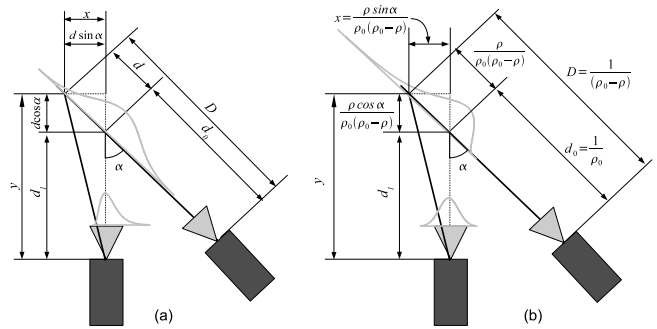


Fig. 3. Uncertainty propagation from the scene point to the image. (a) XYZ coding. (b) Inverse depth coding.

point from a distance d_1 ; the parallax angle α is approximated by the angle between the cameras' optic axes.

The point's location error, d , is encoded as Gaussian in depth:

$$D = d_0 + d, \quad d \sim N(0, \sigma_d^2). \quad (13)$$

This error d is propagated to the image of the point in the second camera, u as:

$$u = \frac{x}{y} = \frac{d \sin \alpha}{d_1 + d \cos \alpha}. \quad (14)$$

The Gaussianity of u is analyzed by means of (12), giving linearity index:

$$L_d = \left| \frac{\frac{\partial^2 u}{\partial d^2} 2\sigma_d}{\frac{\partial u}{\partial d}} \right| = \frac{4\sigma_d}{d_1} |\cos \alpha| \quad (15)$$

C. Linearity of Inverse Depth Parametrization

The inverse depth parametrization is based on the same scene geometry as the direct depth coding, but the depth error is encoded as Gaussian in inverse depth (Fig 3b):

$$D = \frac{1}{\rho_0 - \rho}, \quad \rho \sim N(0, \sigma_\rho^2) \quad (16)$$

$$d = D - d_0 = \frac{\rho}{\rho_0(\rho_0 - \rho)}, \quad d_0 = \frac{1}{\rho_0}. \quad (17)$$

So the image of the scene point is computed as:

$$u = \frac{x}{y} = \frac{d \sin \alpha}{d_1 + d \cos \alpha} = \frac{\rho \sin \alpha}{\rho_0 d_1 (\rho_0 - \rho) + \rho \cos \alpha} \quad (18)$$

and the linearity index L_ρ is now:

$$L_\rho = \left| \frac{\frac{\partial^2 u}{\partial \rho^2} 2\sigma_\rho}{\frac{\partial u}{\partial \rho}} \right| = \frac{4\sigma_\rho}{\rho_0} \left| 1 - \frac{d_0}{d_1} \cos \alpha \right|. \quad (19)$$

D. Depth vs. Inverse Depth Comparison

When a feature is initialized, the depth prior has to cover a vast region in front of the camera. With the inverse depth representation, the 95% confidence region with parameters ρ_0, σ_ρ is:

$$\left[\frac{1}{\rho_0 + 2\sigma_\rho}, \frac{1}{\rho_0 - 2\sigma_\rho} \right]. \quad (20)$$

This region cannot include zero depth but can easily extend to infinity.

Conversely, with the depth representation the 95% region with parameters d_0, σ_d is $[d_0 - 2\sigma_d, d_0 + 2\sigma_d]$. This region can include zero depth but cannot extend to infinity.

In the first few frames after a new feature has been initialized, little parallax is likely to have been observed. Therefore $\frac{d_0}{d_1} \approx 1$ and $\alpha \approx 0 \implies \cos \alpha \approx 1$. In this case the L_d linearity index for depth is high (bad), while the L_ρ linearity index for inverse depth is low (good): during initialization the inverse depth measurement equation linearity is superior to the XYZ coding.

As estimation proceeds and α increases, leading to more accurate depth estimates, the inverse depth representation continues to have a high degree of linearity. This is because in the expression for L_ρ the increase in the term $|1 - \frac{d_0}{d_1} \cos \alpha|$ is compensated by the decrease in $\frac{4\sigma_\rho}{\rho_0}$. For inverse depth features a good linearity index is achieved along the whole estimation history. So the inverse depth coding is suitable for both low and high parallax cases if the feature is continuously observed.

The XYZ encoding has low computational cost, but achieves linearity only at low depth uncertainty and high parallax. In Section VI we explain how the representation of a feature can be switched over such that the inverse depth parametrization is only used when needed — for features which are either just initialized or at extreme depths.

V. FEATURE INITIALIZATION

From just a single observation no feature depth can be estimated (although it would be possible in principle to impose a very weak depth prior by knowledge of the type of scene observed). What we do is to assign a general Gaussian prior in inverse depth which encodes probabilistically the fact that the point has to be in front of the camera. Hence, thanks to the linearity of inverse depth at low parallax, the filter can be initialized from just one observation. Experimental tuning has shown that infinity should be included with reasonable probability within the initialization prior, despite the fact that this means that depth estimates can become negative. Once initialized, features are processed with the standard EKF prediction-update loop — even in the case of negative inverse depth estimates — and immediately contribute to camera location estimation within SLAM.

It is worth noting that while a feature retains low parallax, it will automatically be used mainly to determine the camera orientation. The feature's depth will remain uncertain, with the hypothesis of infinity still under consideration (represented by the probability mass corresponding to negative inverse depths). If the camera translates to produce enough parallax then the feature's depth estimation will be improved and it will begin to contribute more to camera location estimation.

The initial location for a newly observed feature inserted into the state vector is:

$$\hat{\mathbf{y}} \left(\hat{\mathbf{r}}^{WC}, \hat{\mathbf{q}}^{WC}, \mathbf{h}, \rho_0 \right) = \left(\hat{x}_i \quad \hat{y}_i \quad \hat{z}_i \quad \hat{\theta}_i \quad \hat{\phi}_i \quad \hat{\rho}_i \right)^\top, \quad (21)$$

a function of the current camera pose estimate $\hat{\mathbf{r}}^{WC}, \hat{\mathbf{q}}^{WC}$, the image observation $\mathbf{h} = (u \ v)^\top$ and the parameters determining the depth prior ρ_0, σ_ρ .

The end-point of the initialization ray (see Figure 1) is taken from the current camera location estimate:

$$\left(\hat{x}_i \quad \hat{y}_i \quad \hat{z}_i \right)^\top = \hat{\mathbf{r}}^{WC}, \quad (22)$$

and the direction of the ray is computed from the observed point, expressed in the world coordinate frame:

$$\mathbf{h}^W = \mathbf{R}_{WC} \left(\mathbf{q}^{WC} \right) \left(v \quad \nu \quad 1 \right)^\top, \quad (23)$$

where v and ν are normalized retina image coordinates. Despite \mathbf{h}^W being a non-unit directional vector, the angles by which we parametrize its direction can be calculated as:

$$\begin{pmatrix} \theta_i \\ \phi_i \end{pmatrix} = \begin{pmatrix} \arctan \left(\frac{\mathbf{h}_x^W, \mathbf{h}_z^W}{\mathbf{h}_y^W, \sqrt{\mathbf{h}_x^{W2} + \mathbf{h}_z^{W2}}} \right) \\ \arctan \left(-\mathbf{h}_y^W, \sqrt{\mathbf{h}_x^{W2} + \mathbf{h}_z^{W2}} \right) \end{pmatrix} \quad (24)$$

The covariance of $\hat{x}_i, \hat{y}_i, \hat{z}_i, \hat{\theta}_i$ and $\hat{\phi}_i$ is derived from the image measurement error covariance \mathbf{R}_i and the state covariance estimate $\hat{\mathbf{P}}_{k|k}$.

The initial value for ρ_0 and its standard deviation are set empirically such that the 95% confidence region spans a range of depths from close to the camera out to infinity. In our experiments we set $\hat{\rho}_0 = 0.1, \sigma_\rho = 0.5$, which gives an inverse depth confidence region $[1.1, -0.9]$. Notice that infinity is included in this range. Experimental validation has shown that the precise values of these parameters are relatively unimportant to the accurate operation of the filter, as long as infinity is clearly included in the confidence interval.

The state covariance after feature initialization is:

$$\hat{\mathbf{P}}_{k|k}^{\text{new}} = \mathbf{J} \begin{pmatrix} \hat{\mathbf{P}}_{k|k} & 0 & 0 \\ 0 & \mathbf{R}_i & 0 \\ 0 & 0 & \sigma_\rho^2 \end{pmatrix} \mathbf{J}^\top \quad (25)$$

$$\mathbf{J} = \left(\begin{array}{c|c} I & 0 \\ \hline \frac{\partial \mathbf{y}}{\partial \mathbf{r}^{WC}}, \frac{\partial \mathbf{y}}{\partial \mathbf{q}^{WC}}, 0, \dots, 0 & \frac{\partial \mathbf{y}}{\partial \mathbf{h}}, \frac{\partial \mathbf{y}}{\partial \rho} \end{array} \right). \quad (26)$$

The inherent scale ambiguity in monocular SLAM has usually been fixed by observing some known initial features that fix the scale (e.g. [8]). A very interesting experimental observation we have made using the inverse depth scheme is that sequential monocular SLAM can operate successfully without *any* known features in the scene, and in fact the experiments we present in this paper do not use an initialization target. In this case, of course the overall scale of the reconstruction and camera motion is undetermined, although with the formulation of the current paper the estimation will settle on a (meaningless) scale of some value. In very recent work [6] we have investigated this issue with a new dimensionless formulation of monocular SLAM.

VI. SWITCHING FROM INVERSE DEPTH TO XYZ

While the inverse depth encoding can be used at both low and high parallax, it is advantageous for reasons of computational efficiency to restrict inverse depth to cases where the XYZ encoding exhibits non linearity according to the L_d index. This section details switching from inverse depth to XYZ for high parallax features.

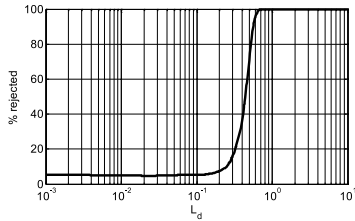


Fig. 4. Percentage of test rejections as a function of the linearity index L_d

A. Conversion from Inverse Depth to XYZ Coding

After each estimation step, the linearity index L_d (Equation 15) is computed for every map feature coded in inverse depth:

$$\mathbf{h}_{\text{XYZ}}^W = \hat{\mathbf{x}}_i - \hat{\mathbf{r}}^{WC}, \quad \sigma_d = \frac{\sigma_\rho}{\rho_i^2}, \quad \sigma_\rho = \sqrt{\mathbf{P}_{\mathbf{y}_i \mathbf{y}_i} (6, 6)}$$

$$d_i = \left\| \mathbf{h}_{\text{XYZ}}^W \right\|, \quad \cos \alpha = \mathbf{m}^\top \mathbf{h}_{\text{XYZ}}^W \left\| \mathbf{h}_{\text{XYZ}}^W \right\|^{-1}. \quad (27)$$

where $\hat{\mathbf{x}}_i$ is computed using equation (4) and $\mathbf{P}_{\mathbf{y}_i \mathbf{y}_i}$ is the submatrix 6×6 covariance matrix corresponding the considered feature.

If L_d is below a switching threshold, the feature is transformed using Equation (4) and the *full state* covariance matrix \mathbf{P} is transformed with the corresponding Jacobian:

$$\mathbf{P}_{\text{new}} = \mathbf{J} \mathbf{P} \mathbf{J}^\top, \quad \mathbf{J} = \text{diag} \left(\mathbf{I}, \frac{\partial \mathbf{x}_i}{\partial \mathbf{y}_i}, \mathbf{I} \right). \quad (28)$$

B. Linearity Index Threshold

We propose to use index L_d (15) to define a threshold for switching from inverse depth to XYZ encoding at the point when the latter can be considered linear. If the XYZ representation is linear, then the measurement u is Gaussian distributed (Equation 10):

$$u \sim N(\mu_u, \sigma_u^2), \quad \mu_u = 0, \quad \sigma_u^2 = \left(\frac{\sin \alpha}{d_1} \right)^2 \sigma_d^2. \quad (29)$$

To determine the threshold in L_d which signals a lack of linearity in the measurement equation a simulation experiment has been performed. The goal was to generate samples from the uncertain distribution for variable u and then apply a standard Kolmogorov-Smirnov Gaussianity [4] test to these samples, counting the percentage of rejected hypotheses, h . When u is effectively Gaussian, the percentage should match the test significance level α_{sl} (5% in our experiments); as the number of rejected hypotheses increases the measurement equation departs from linearity. A plot of the percentage of rejected hypotheses h with respect to the linearity index L_d is shown in Figure 4. It can be clearly seen that when $L_d > 0.2$, h sharply departs from 5%. So we propose the $L_d < 10\%$ threshold for switching from inverse depth to XYZ encoding.

Notice that the plot in Figure 4 is smooth (log scale in L_d), which indicates that the linearity index effectively represents the departure from linearity.

The simulation has been performed for a variety of values of α , d_1 and σ_d ; more precisely all triplets resulting from the following parameter values:

$$\begin{aligned} \alpha(\text{deg}) &\in \{0.1, 1, 3, 5, 7, 10, 20, 30, 40, 50, 60, 70\} \\ d_1(\text{m}) &\in \{1, 3, 5, 7, 10, 20, 50, 100\} \\ \sigma_d(\text{m}) &\in \{0.05, 0.1, 0.25, 0.5, 0.75, 1, 2, 5\}. \end{aligned}$$

The simulation algorithm detailed in Figure 5 is applied to every triplet $\{\alpha, d_1, \sigma_d\}$ to count the percentage of rejected hypotheses h and the corresponding linearity index L_d .

```

input:  $\alpha, d_1, \sigma_d$ 
output:  $h, L_d$ 
 $\sigma_u = \left| \frac{\sin \alpha}{d_1} \right| \sigma_d; \mu_u = 0; // (29)$ 
 $\alpha_{sl} = 0.05; // \text{Kolm. test sign. level}$ 
 $L_d = \frac{4\sigma_d}{d_1} |\cos \alpha|$ 
n_rejected=0;
N_GENERATED_SAMPLES=1000;
SAMPLE_SIZE=1000;

for j=1 to N_GENERATED_SAMPLES repeat
   $\{d_i\}_j = \text{random\_normal}(0, \sigma_d^2, \text{SAMPLE\_SIZE});$ 
  //generate a normal sample from  $N(0, \sigma_d^2)$ ;
   $\{u_i\}_j = \text{propagate\_from\_dept\_to\_image}(\{d_i\}_j, \alpha, d_1); // (14)$ 
  if rejected==KolmogorovSmirnov( $\{u_i\}_j, \mu_u, \sigma_u, \alpha_{sl}$ )
    n_rejected=n_rejected+1;
endfor
 $h = 100 \frac{[n\_rejected]}{[N\_GENERATED\_SAMPLES]};$ 

```

Fig. 5. Simulation algorithm to test the linearity of the measurement equation.

VII. EXPERIMENTAL RESULTS

The performance of the new parametrization has been tested on real image sequences acquired with a hand-held low cost Unibrain IEEE1394 camera, with a 90° field of view and 320×240 resolution, capturing monochrome image sequences at 30 fps.

Five experiments were performed. The first was an indoor sequence processed offline with a Matlab implementation, the goal being to analyze initialization of scene features located at different depths. The second experiment shows an outdoor sequence processed in real-time with a C++ implementation. The focus was on distant features, observed under low parallax along the whole sequence. The third experiment was a loop closing sequence, concentrating on camera covariance evolution. Fourth was a simulation experiment to analyze the effect of switching from inverse depth to XYZ representations. In the last experiment the switching performance was verified on the real loop closing sequence. This section ends with a computing time analysis. It is worth noting that no initial pattern to fix the scale was used in any of the last three experiments.

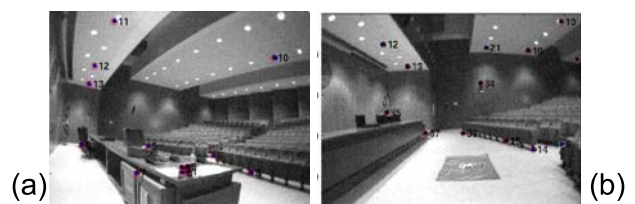


Fig. 6. First (a) and last (b) frame in the sequence of the indoor experiment of Section VII-A. Features 11, 12, 13 are analyzed. These features are initialized in the same frame but are located at different distances from the camera.

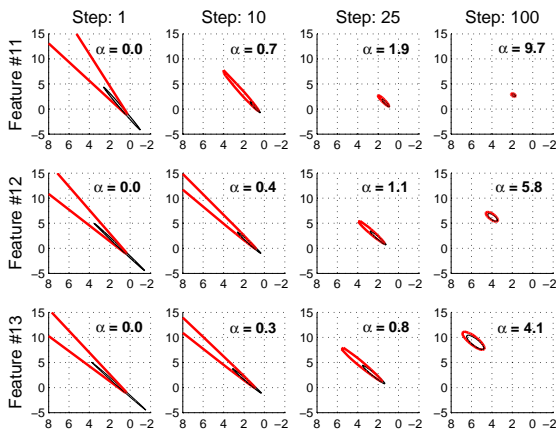


Fig. 7. Feature initialization. Each column shows the estimation history for a feature horizontal components. For each feature, the estimates after 1, 10, 25 and 100 frames since initialization are plotted; the parallax angle α in degrees between the initial observation and the current frame is displayed. The thick (red) lines show (calculated by a Monte Carlo numerical simulation) the 95% confidence region when coded as Gaussian in inverse depth. The thin (black) ellipsoids show the uncertainty as a Gaussian in XYZ space propagated according to Equation (28). Notice how at low parallax the inverse depth confidence region is very different from the elliptical Gaussian. However, as the parallax increases, the uncertainty reduces and collapses to the Gaussian ellipse.

A. Indoor Sequence

This experiment analyzes the performance of the inverse depth scheme as several features at a range of depths are tracked within SLAM. We discuss three features, which are all detected in the same frame but have very different depths. Figure 6 shows the image where the analyzed features are initialized (frame 18 in the sequence) and the last image in the sequence. Figure 7 focuses on the evolution of the estimates corresponding to the features, with labels 11, 12 and 13, at frames 1, 10, 25 and 100. Confidence regions derived from the inverse depth representation (thick red line) are plotted in XYZ space by numerical Monte Carlo propagation from the six-dimensional multivariate Gaussians representing these features in the SLAM EKF. For comparison, standard Gaussian XYZ acceptance ellipsoids (thin black line) are linearly propagated from the six-dimensional representation by means of the Jacobian of equation (28). The parallax α in degrees for each feature at every step is also displayed.

When initialized, the 95% acceptance region of all the features includes $\rho = 0$ so infinite depth is considered as a possibility. The corresponding confidence region in depth is highly asymmetric, excluding low depths but extending to infinity. It is clear that Gaussianity in inverse depth is not mapped to Gaussianity in XYZ, so the black ellipsoids produced by Jacobian transformation are far from representing the true depth uncertainty. As stated in Section IV-D, it is at low parallax that the inverse depth parametrization plays a key role.

As rays producing bigger parallax are gathered, the uncertainty in ρ becomes smaller but still maps to a non-Gaussian distribution in XYZ. Eventually, at high parallax, for all of the features the red confidence regions become closely Gaussian and well-approximated by the linearly-propagated

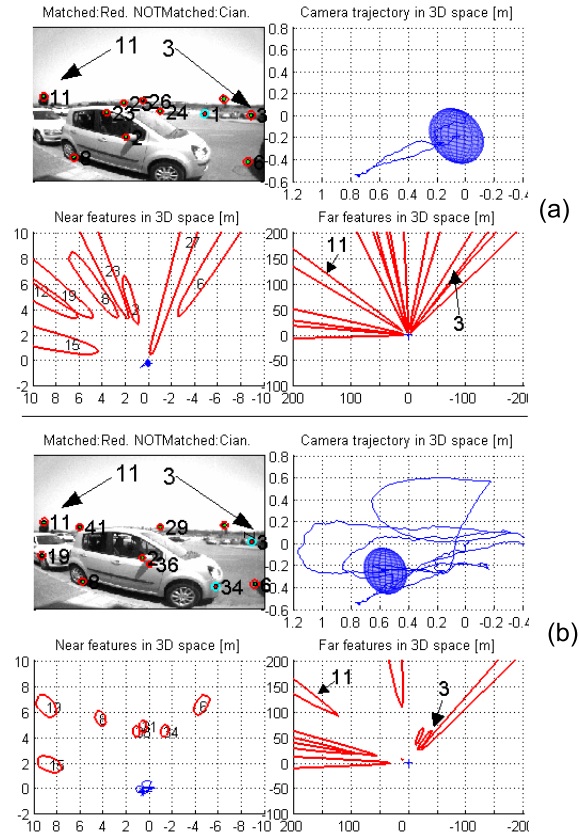


Fig. 8. Subfigures (a) and (b) show frames #163 and #807 from the outdoor experiment of Section VII-B. This experiment was processed in real time. The focus was two features: 11 (tree on the left) and 3 (car on the right) at low parallax. Each of the two subfigures shows the current images, and top-down views illustrating the horizontal components of the estimation of camera and feature locations at three different zoom scales for clarity: the top-right plots (maximum zoom) highlight the estimation of the camera motion; bottom-left (medium zoom) views highlight nearby features; and bottom-right (minimum zoom) emphasizes distant features.

black ellipses — but this happens much sooner for nearby feature 11 than distant feature 13.

A movie showing the input sequence and estimation history of this experiment is available as multimedia data `inverseDepth_indoor.avi`. The raw input image sequence is also available as `inverseDepth_indoorRawImages.tar.gz`.

B. Real-Time Outdoor Sequence

This 860 frame experiment was performed with a C++ implementation which achieves real-time performance at 30 fps with hand-held camera. Here we highlight the ability of our parametrization to deal with both close and distant features in an outdoor setting. The input image sequence is available as multimedia material `inverseDepth_outdoorRawImages.tar.gz`. A movie showing the estimation process is also available as `inverseDepth_outdoor.avi`.

Figure 8 shows two frames of the movie illustrating the performance. For most of the features, the camera ended up gathering enough parallax to accurately estimate their depths. However, being outdoors, there were distant features producing low parallax during the whole camera motion.

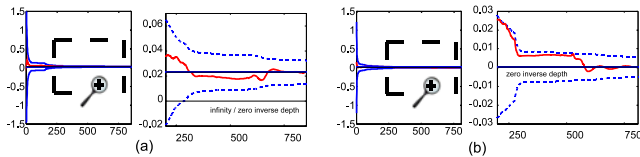


Fig. 9. Analysis of outdoor experiment of Section VII-B. (a) Inverse depth estimation history for feature 3, on the car, and (b) for feature 11, on a distant tree. Due to the uncertainty reduction during estimation, two plots at different scales are shown for each feature. It is shown the 95% confidence region, and with a thick line the estimated inverse depth. The thin solid line is the inverse depth estimated after processing the whole sequence. In (a), for the first 250 steps, zero inverse depth is included in confidence region, meaning that the feature might be at infinity. After this, more distant but finite locations are gradually eliminated, and eventually the feature's depth is accurately estimated. In (b), the tree is so distant that the confidence region always includes zero, since little parallax is gathered for that feature.

The inverse depth estimation history for two features is highlighted in Figure 9. It is shown that distant, low parallax features are persistently tracked through the sequence, despite the fact that their depths cannot be precisely estimated. The large depth uncertainty, represented with the inverse depth scheme, is successfully managed by the SLAM EKF, allowing the orientation information supplied by these features to be exploited.

Feature 3, on a nearby car, eventually gathers enough parallax enough to have an accurate depth estimate after 250 images where infinite depth is considered as a possibility. Meanwhile the estimate of Feature 11, on a distant tree and never displaying significant parallax, never collapses in this way and zero inverse depth remains within its confidence region. Delayed initialization schemes would have discarded this feature without obtaining any information from it, while in our system it behaves like a bearing reference. This ability to deal with distant points in real time is a highly advantageous quality of our parametrization. Note that what does happen to the estimate of Feature 11 as translation occurs is that hypotheses of nearby depths are ruled out — the inverse depth scheme correctly recognizes that measuring little parallax while the camera has translated some distance allows a minimum depth for the feature to be set.

C. Loop Closing Sequence

A loop closing sequence offers a challenging benchmark for any SLAM algorithm. In this experiment a handheld camera was carried by a person walking in small circles within a very large student laboratory, carrying out two complete laps. The raw input image sequence is available as `inverseDepth_loopClosingRawImages.tar.gz`, and a movie showing the mapping process as `inverseDepth_loopClosing.avi`.

Figure 10 shows a selection of the 737 frames from the sequence, concentrating on the beginning, first loop closure and end of the sequence. Figure 11 shows the camera location estimate covariance history, represented by the 95% confidence regions for the 6 camera d.o.f. and expressed in a reference local to the camera.

We observe the following properties of the evolution of the estimation, focussing in particular on the uncertainty in the camera location:

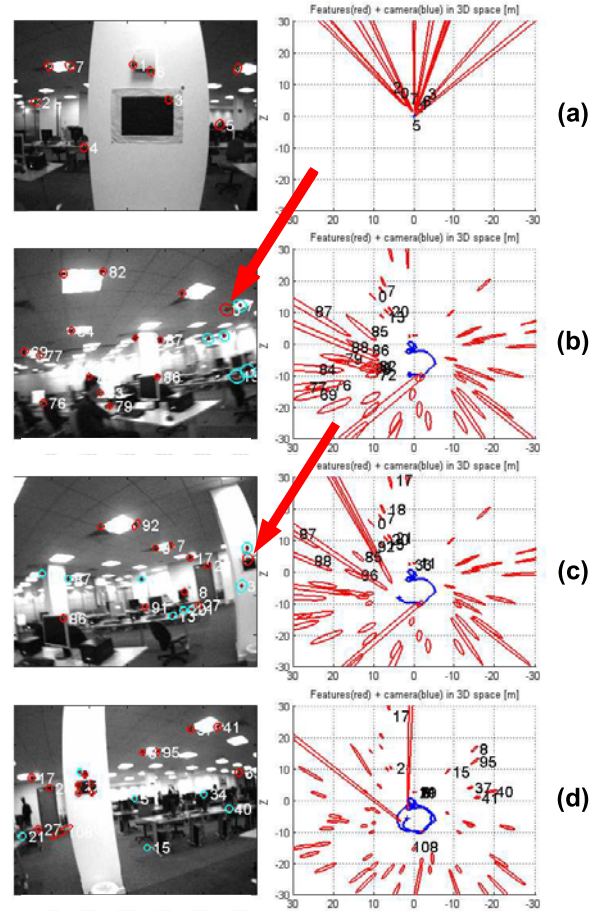


Fig. 10. A selection of frames from the loop closing experiment of Section VII-C. For each frame, we show the current image and reprojected map (left), and a top-down view of the map with 95% confidence regions and camera trajectory (right). Notice that confidence regions for the map features are far from being Gaussian ellipses, especially for newly initialized or distant features. The selected frames are: (a) #11, close to the start of the sequence; (b) #417, where the first loop closing match, corresponding to a distant feature, is detected; the loop closing match is signaled with an arrow; (c) #441 where the first loop closing match corresponding to a close feature is detected; the match is signaled with an arrow; and (d) #737, the last image, in the sequence, after reobserving most of the map features during the second lap around the loop.

- After processing the first few images, the uncertainty in the depth of features is huge, with highly non-elliptical confidence regions in XYZ space (Fig. 10(a)).
- In Figure 11 the first peak in the X and Z translation uncertainty corresponds to a camera motion backwards along the optical axis; this motion produces poor parallax for newly initialized features, and we therefore see a reduction in orientation uncertainty and an increase in translation uncertainty. After frame #50 the camera again translates in the X direction, parallax is gathered and the translation uncertainty is reduced.
- From frame #240, the camera starts a 360° circular motion in the XZ plane. The camera explores new scene regions, and the covariance increases steadily as expected (Fig. 11).
- In frame #417, the first loop closing feature is re-observed. This is a feature which is distant from the camera, and causes an abrupt reduction in orientation and

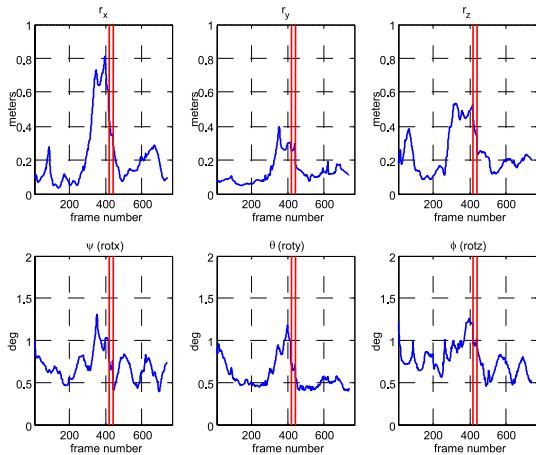


Fig. 11. Camera location estimate covariance along the sequence. The 95% confidence regions for each of the 6 d.o.f of camera motion are plotted. Note that errors are expressed in a reference local to the camera. The vertical solid lines indicate the loop closing frames #417 and #441.

translation uncertainty (Fig. 11), though a medium level of uncertainty remains.

- In frame #441, a much closer loop closing feature (mapped with high parallax) is matched. Another abrupt covariance reduction takes place (Fig. 11) with the extra information this provides.
- After frame #441, as the camera goes on a second lap around the loop, most of the map features are revisited, almost no new features are initialized, and hence the uncertainty in the map is further reduced. Comparing the map at frame #441 (the beginning of the second lap) and at #737, (the end of the second lap), we see a significant reduction in uncertainty. During the second lap, the camera uncertainty is low, and as features are reobserved their uncertainties are noticeably reduced (Fig. 10(c) and (d)).

Note that these loop closing results with the inverse depth representation show a marked improvement on the experiments on monocular SLAM with a humanoid robot presented in [9], where a gyro was needed in order to reduce angular uncertainty enough to close loops with very similar camera motions.

D. Simulation Analysis for Inverse Depth to XYZ Switching

In order to analyze the effect of the parametrization switching proposed in Section VI on the consistency of SLAM estimation, simulation experiments with different switching thresholds were run. In the simulations, a camera completed two laps of a circular trajectory of radius 3m in the XZ plane, looking out radially at a scene composed of points lying on three concentric spheres of radius 4.3m, 10m and 20m. These points at different depths were intended to produce observations with a range of parallax angles (Figure 12.)

The camera parameters of the simulation correspond with our real image acquisition system: camera 320×240 pixels,

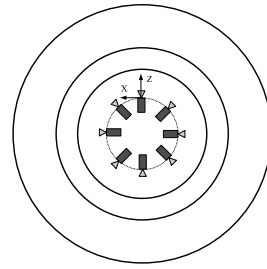


Fig. 12. Simulation configuration for analysis of parametrization switching in Section VII-D, sketching the circular camera trajectory and 3D scene, composed of three concentric spheres of radius 4.3m, 10m and 20m. The camera completes two circular laps in the (XZ) plane with radius 3m, and is orientated radially.

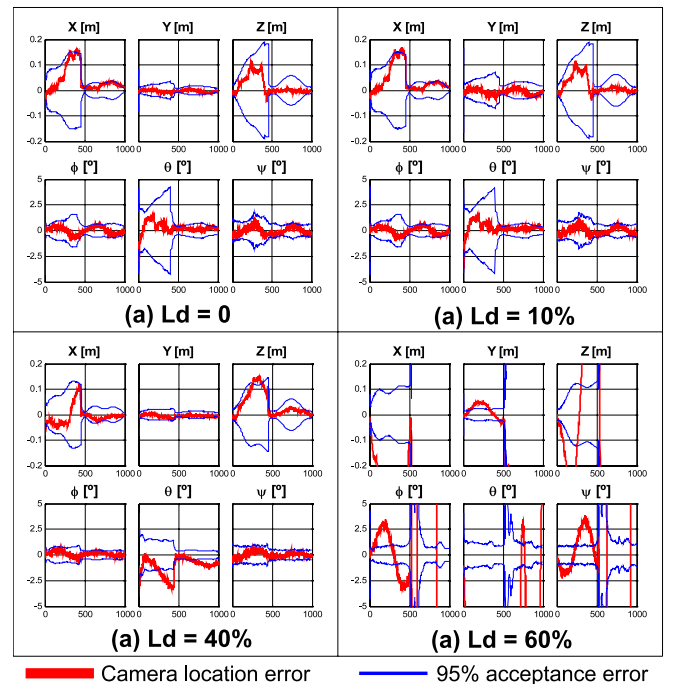


Fig. 13. Details from the parametrization switching experiment. Camera location estimation error history in 6 d.o.f. (translation in XYZ , and three orientation angles $\psi\theta\phi$) for four switching thresholds: With $L_d = 0\%$, no switching occurs and the features all remain in the inverse depth parametrization. At, $L_d = 10\%$ although features from the spheres at 4.3m and 10m are eventually converted, no degradation with respect to the non-switching case is observed. At $L_d = 40\%$ some features are switched before achieving true Gaussian, and there is noticeable degradation, especially in θ rotation around the Y axis. At $L_d = 60\%$ the map becomes totally inconsistent and loop closing fails.

frame rate 30 frames/sec, image field of view 90° , measurement uncertainty for a point feature in the image, Gaussian $N(0, 1\text{pixel}^2)$. The simulated image sequence contained 1000 frames. Features were selected following the randomized map management algorithm proposed in [8] in order to have 15 features visible in the image at all times. All our simulation experiments work using the same scene features, in order to homogenize the comparison.

Four simulation experiments for different thresholds for switching each feature from inverse depth to XYZ parametrization were run, with $L_d \in \{0\%, 10\%, 40\%, 60\%\}$. Figure 13 shows the camera trajectory estimation history in

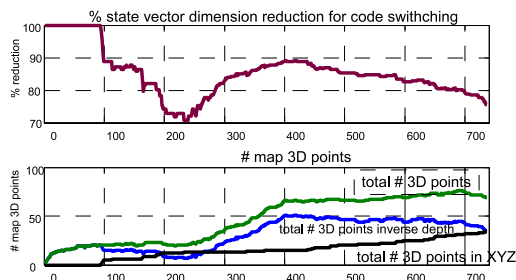


Fig. 14. Parametrization switching on a real sequence (Section VII-E): state vector size history. Top: percentage reduction in state dimension when using switching compared with keeping all points in inverse depth. Bottom: total number of points in the map, showing the number of points in inverse depth and the number of points in XYZ.

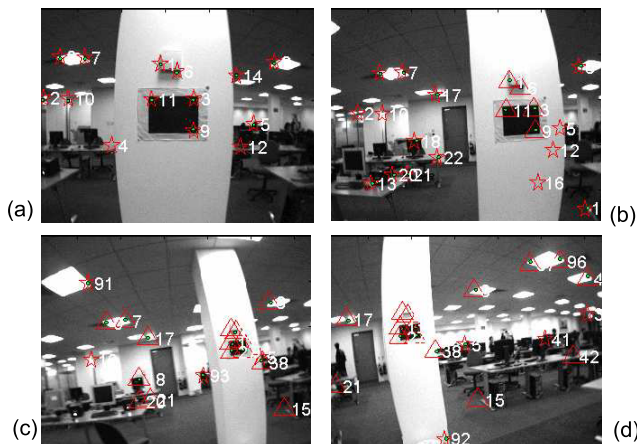


Fig. 15. Parametrization switching seen in image space: points coded in inverse depth (\star) and coded in XYZ (Δ). (a) First frame, with all features in inverse depth. (b) Frame #100; nearby features start switching. (c) Frame # 470, loop closing; most features in XYZ. (d) Last image of the sequence.

6 d.o.f. (translation in XYZ , and three orientation angles $\psi(\text{Rot}_x)$, $\theta(\text{Rot}_y)$, $\phi(\text{Rot}_z)$, cyclotorsion)). The following conclusions can be made:

- Almost the same performance is achieved with no switching (0%) and with 10% switching. So it is clearly advantageous to perform 10% switching because there is no penalization in accuracy and the state vector size of each converted feature is halved.
- Switching too early degrades accuracy, especially in the orientation estimate. Notice how for 40% the orientation estimate is worse and the orientation error covariance is smaller, showing filter inconsistency. For 60%, the estimation is totally inconsistent and loop closing fails.
- Since early switching degrades performance, the inverse depth parametrization is mandatory for initialization of every feature and over the long-term for low-parallax features.

E. Parametrization Switching with Real Images

The loop closing sequence of Section VII-C was processed without any parametrization switching, and with switching at $L_d = 10\%$. A movie showing the results is available as `inverseDepth_loopClosing_ID_to_XYZ_conversion.avi`.

As in the simulation experiments, no significant change was noticed in the estimated trajectory or map.

Figure 14 shows the history of the state size, the number of map features and how their parametrization evolves. At the last estimation step about half of the features had been switched; at this step the state size had reduced from 427 to 322 (34 inverse depth features and 35 XYZ), i.e. 75% of the original vector size. Figure 15 shows four frames from the sequence illustrating feature switching. Up to step 100 the camera has low translation and all the features are in inverse depth form. As the camera translates nearby features switch to XYZ. Around step 420, the loop is closed and features are reobserved, producing a significant reduction in uncertainty which allows switching of more reobserved close features. Our method automatically determines which features should be represented in the inverse depth or XYZ forms, optimizing computational efficiency without sacrificing accuracy.

F. Processing Time

We give some details of the real-time operation of our monocular SLAM system, running on a 1.8 GHz. Pentium M processor laptop. A typical EKF iteration would imply:

- A state vector dimension of 300.
- 12 features observed in the image, a measurement dimension of 24.
- 30 fps, so 33.3 ms available for processing.

Typical computing time breaks down as follows: Image acquisition, 1 ms.; EKF prediction, 2 ms.; Image matching, 1 ms.; EKF update, 17 ms. That adds up to a total of 21ms. The remaining time is used for graphics functions, using OpenGL on an Nvidia card and scheduled at a low priority.

The quoted state vector size 300 corresponds to a map size of 50 if all features are encoded using inverse depth. In indoor scenes, thanks to switching maps of up to 60-70 features can be computed in real time. This size is enough to map many typical scenes robustly.

VIII. CONCLUSION

We have presented a parametrization for monocular SLAM which permits operation based uniquely on the standard EKF prediction-update procedure at every step, unifying initialization with the tracking of mapped features. Our inverse depth parametrization for 3D points allows unified modelling and processing for any point in the scene, close or distant, or even at ‘infinity’. In fact, close, distant or just-initialized features are processed within the routine EKF prediction-update loop without making any binary decisions. Thanks to the undelayed initialization and immediate full use of infinite points, estimates of camera orientation are significantly improved, reducing the camera estimation jitter often reported in previous work. The jitter reduction in turn leads to computational benefits in terms of smaller search regions and improved image processing speed

The key factor is that due to our parametrization of the direction and inverse depth of a point relative to the location from which it was first seen, our measurement equation has

low linearization errors at low parallax, and hence the estimation uncertainty is accurately modeled with a multi-variate Gaussian. In Section IV we presented a model which quantifies linearization error. This provides a theoretical understanding of the impressive outdoor, real-time performance of the EKF with our parametrization.

The inverse depth representation requires a six-dimensional state vector per feature, compared to three for XYZ coding. This doubles the map state vector size, and hence produces a 4-fold increase in the computational cost of the EKF. Our experiments show that it is essential to retain the inverse depth parametrization for initialization and distant features, but that nearby features can be safely converted to the cheaper XYZ representation meaning that the long-term computational cost need not significantly increase. We have given details on when this conversion should be carried out for each feature, to optimize computational efficiency without sacrificing accuracy.

The experiments presented have validated the method with real imagery, using a hand-held camera as the only sensor both indoors and outdoors. We have experimentally verified the key contributions of our work:

- Real-time performance achieving 30 fps real-time processing for maps up to 60–70 features.
- Real-time loop closing.
- Dealing simultaneously with low and high parallax features.
- Non delayed initialization.
- Low jitter, full 6 DOF monocular SLAM.

In the experiments, we have focused on a map size around 60–100 features, because these map sizes can be dealt with in real time at 30Hz and we have focused on the challenging loop closing issue. Useful future work would be a thorough analysis of the the limiting factors in EKF inverse depth monocular SLAM in terms of linearity, data association errors, accuracy, map size and ability to deal with degenerate motion such as pure rotations or a static camera for long time periods.

Finally, our simulations and experiments have shown that inverse depth monocular SLAM operates well without known patterns in the scene to fix scale. This result points towards further work in understanding the role of scale in monocular SLAM (an avenue we have begun to investigate in a dimensionless formulation in [6]) and in further bridging the gap between sequential SLAM techniques and structure from motion methods from the computer vision literature.

APPENDIX

To recover the ideal projective undistorted coordinates $\mathbf{h}_u = (u_u, v_u)^\top$, from the actually distorted ones gathered by the camera, $\mathbf{h}_d = (u_d, v_d)^\top$, the classical two parameters radial distortion model [17] is applied:

$$\begin{pmatrix} u_d \\ v_d \end{pmatrix} = \mathbf{h}_u \begin{pmatrix} u_d \\ v_d \end{pmatrix} = \begin{pmatrix} u_0 + (u_d - u_0) (1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) \\ v_0 + (v_d - v_0) (1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) \end{pmatrix} \quad (30)$$

$$r_d = \sqrt{(d_x (u_d - u_0))^2 + (d_y (v_d - v_0))^2}$$

Where, u_0, v_0 are the image center and, κ_1, κ_2 are the radial distortion coefficients.

To compute the distorted coordinates from the undistorted:

$$\begin{pmatrix} u_d \\ v_d \end{pmatrix} = \mathbf{h}_d \begin{pmatrix} u_u \\ v_u \end{pmatrix} = \begin{pmatrix} u_0 + \frac{(u_u - u_0)}{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4)} \\ v_0 + \frac{(v_u - v_0)}{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4)} \end{pmatrix} \quad (31)$$

$$r_u = r_d (1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) \quad (32)$$

$$r_u = \sqrt{(d_x (u_u - u_0))^2 + (d_y (v_u - v_0))^2} \quad (33)$$

r_u is readily computed from (33), but r_d has to be numerically solved from (32), e.g using Newton-Raphson, hence (31) can be used to compute the distorted point.

Undistortion jacobian, $\frac{\partial \mathbf{h}_u}{\partial (u_d, v_d)}$ has analytical expression:

$$\begin{pmatrix} \frac{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) + 2((u_d - u_0) d_x)^2 \times (\kappa_1 + 2\kappa_2 r_d^2)}{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) + 2((v_d - v_0) d_y)^2 \times (\kappa_1 + 2\kappa_2 r_d^2)} & \frac{2d_y^2 (u_d - u_0) (v_d - v_0) \times (\kappa_1 + 2\kappa_2 r_d^2)}{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) + 2((v_d - v_0) d_y)^2 \times (\kappa_1 + 2\kappa_2 r_d^2)} \\ \frac{2d_x^2 (v_d - v_0) (u_d - u_0) \times (\kappa_1 + 2\kappa_2 r_d^2)}{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) + 2((u_d - u_0) d_x)^2 \times (\kappa_1 + 2\kappa_2 r_d^2)} & \frac{2d_x^2 (v_d - v_0) (u_d - u_0) \times (\kappa_1 + 2\kappa_2 r_d^2)}{(1 + \kappa_1 r_d^2 + \kappa_2 r_d^4) + 2((u_d - u_0) d_x)^2 \times (\kappa_1 + 2\kappa_2 r_d^2)} \end{pmatrix} \quad (34)$$

The jacobian for the distortion is computed by inverting expression (34):

$$\frac{\partial \mathbf{h}_d}{\partial (u_{u_0}, v_{u_0})} \Big|_{(u_u, v_u)} = \left(\frac{\partial \mathbf{h}_u}{\partial (u_d, v_d)} \Big|_{\mathbf{h}_d(u_{u_0}, v_{u_0})} \right)^{-1} \quad (35)$$

ACKNOWLEDGMENT

We are very grateful to David Murray, Ian Reid and other members of Oxford's Active Vision Laboratory for discussions and software collaboration.

Thank you to the anonymous reviewers for their useful comments.

REFERENCES

- [1] V. J. Aidala and S. E. Hammel. Utilization of modified polar coordinates for bearing-only tracking. *IEEE Trans. Autom. Control*, 28(3):283–294, March 1983.
- [2] T. Bailey. Constrained initialisation for bearing-only SLAM. In *Proc. IEEE Int. Conf. Robotics and Automation*, Taiwan, 2003.
- [3] M. Bryson and S. Sukkarieh. Bearing-only SLAM for an airborne vehicle. In *Australian Conference on Robotics and Automation (ACRA '05)*, Sidney, 2005.
- [4] G. C. Canavos. *Applied Probability and Statistical Methods*. Little, Brown and Company, Boston, USA, 1984.
- [5] A. Chowdhury and R. Chellappa. Stochastic approximation and rate-distortion analysis for robust structure and motion estimation. *International Journal of Computer Vision*, 55(1):27–53, 2003.
- [6] J. Civera, A. J. Davison, and J. M. M. Montiel. Dimensionless monocular SLAM. In *3rd Iberian Conf. on Pattern Recognition and Image Analysis*, 2007.
- [7] J. Civera, A. J. Davison, and J. M. M. Montiel. Inverse depth to depth conversion for monocular SLAM. In *Proc. Intl. Conf. on Robotics and Automation*, pages 2778–2783, 2007.
- [8] A. Davison. Real-time simultaneous localization and mapping with a single camera. In *Proc. International Conference on Computer Vision*, 2003.
- [9] A. J. Davison, I. Reid, N. Molton, and O. Stasse. Real-time single camera SLAM. *IEEE Trans. on PAMI*, 2007.
- [10] E. Eade and T. Drummond. Scalable monocular SLAM. In *In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2006.
- [11] A. W. Fitzgibbon and A. Zisserman. Automatic camera recovery for closed or open image sequences. In *European Conference on Computer Vision*, pages 311–326, 1998.
- [12] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004.

- [13] D. Heeger and A. Jepson. Subspace methods for recovering rigid motion I: Algorithm and implementation. *International Journal of Computer Vision*, pages 95–117, 1992.
- [14] J. H. Kim and S. Sukkarieh. Airborne simultaneous localisation and map building. In *Proceedings of the IEEE International Conference on Robotics and Automation*, pages 406–411, 2003.
- [15] N. Kwok and G. Dissanayake. An efficient multiple hypothesis filter for bearing-only SLAM. In *IROS*, pages 736–741, 2004.
- [16] L. Matthies, T. Kanade, and R. Szeliski. Kalman filter-based algorithms for estimating depth from image sequences. *International Journal of Computer Vision*, 3(3):209–238, 1989.
- [17] E. Mikhail, J. Bethel, and M. J.C. *Introduction to Modern Photogrammetry*. John Wiley & Sons, 2001.
- [18] J. Montiel, J. Civera, and A. J. Davison. Unified inverse depth parametrization for monocular SLAM. In *Robotics Science and Systems Conference*. Philadelphia., 2006.
- [19] J. Montiel and A. J. Davison. A visual compass based on SLAM. In *Proc. Intl. Conf. on Robotics and Automation*, pages 1917–1922, 2006.
- [20] E. Mouragnon, M. Lhuillier, M. Dhome, F. Dekeyser, and P. Sayd. Real-time localization and 3D reconstruction. In *IEEE Conf. on Computer Vision and Pattern Recognition*, pages 1027–1031, 2006.
- [21] D. Nister, O. Naroditsky, and J. Bergen. Visual odometry for ground vehicle applications. *Journal of Field Robotics*, 23(1):3–26, 2006.
- [22] M. Okutomi and T. Kanade. A multiple-baseline stereo. *IEEE Trans. Pattern Anal. Mach. Intell.*, 15(4):353–363, 1993.
- [23] M. Pollefeys, R. Koch, and L. Van Gool. Self-calibration and metric reconstruction inspite of varying and unknown intrinsic camera parameters. *International Journal of Computer Vision*, 32(1):7–25, 1999.
- [24] J. Sola. *Towards Visual Localization, Mapping and Moving Objects Tracking by a Mobile Robot: a Geometric and Probabilistic Approach*. PhD thesis, LAAS-CNRS, 2007.
- [25] J. Sola, A. Monin, M. Devy, and T. Lemaire. Undelayed initialization in bearing only SLAM. In *2005 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2005.
- [26] N. Trawny and S. I. Roumeliotis. A unified framework for nearby and distant landmarks in bearing-only SLAM. In *Proc. Intl. Conf. on Robotics and Automation*, pages 1923–1929, 2006.



Javier Civera was born in Zaragoza, Spain, in 1980. He received the M.S. degree in industrial-electrical engineering from the University in Zaragoza in 2004. He is currently an assistant lecturer at the University of Zaragoza, where he teaches courses in automatic control theory. He is working towards the PhD degree at the Robotics, Perception and Real-Time Group in this University. His research interests include computer vision and mobile robotics.



Andrew J. Davison read physics at the University of Oxford, receiving the BA degree in 1994. In his doctoral research in Oxfords Robotics Research Group under the supervision of Professor David Murray, he developed one of the first robot SLAM systems using vision. On receiving the DPhil degree in 1998, he took up an EU Science and Technology Fellowship and spent two years at AIST in Japan, expanding his work on visual robot navigation. He returned to further postdoctoral work with Dr. Ian Reid at Oxford in 2000, was awarded a five year EPSRC Advanced Research Fellowship in 2002, and moved to Imperial College London in 2005 to take up a lectureship. He continues to work on advancing the basic technology of real-time localization and mapping using vision while collaborating to apply these techniques in robotics and related areas, and has just been awarded a European Research Council Starting Grant.



José M. Martínez Montiel born in Arnedo, Spain. He received the M.S. and Ph.D. degrees in Electrical Engineering from the University of Zaragoza, Spain, in 1991 and 1996, respectively. He is currently an Associate Professor at the University of Zaragoza in charge of Perception and Computer Vision courses. He is member of the the Robotics, Perception and Real-Time Group. He has been awarded Spanish MEC grants to fund research at University of Oxford and Imperial College London. His current interests are computer vision, real time vision localization and mapping research and the transference of this technology to robotic and non robotic application domains.