

Fusing Monocular Information in Multicamera SLAM

Joan Solà, André Monin, Michel Devy, and Teresa Vidal-Calleja

Abstract—This paper explores the possibilities of using monocular simultaneous localization and mapping (SLAM) algorithms in systems with more than one camera. The idea is to combine in a single system the advantages of both monocular vision (bearings-only, infinite range observations but no 3-D instantaneous information) and stereovision (3-D information up to a limited range). Such a system should be able to instantaneously map nearby objects while still considering the bearing information provided by the observation of remote ones. We do this by considering each camera as an independent sensor rather than the entire set as a monolithic supersensor. The visual data are treated by monocular methods and fused by the SLAM filter. Several advantages naturally arise as interesting possibilities, such as the desynchronization of the firing of the sensors, the use of several unequal cameras, self-calibration, and cooperative SLAM with several independently moving cameras. We validate the approach with two different applications: a stereovision SLAM system with automatic self-calibration of the rig's main extrinsic parameters and a cooperative SLAM system with two independent free-moving cameras in an outdoor setting.

Index Terms—Calibration, image sequence analysis, Kalman filtering, machine vision, robot vision systems, stereovision.

I. INTRODUCTION

THE SIMULTANEOUS localization and mapping (SLAM) problem, as formulated by the robotics community, is that of creating a *map* of the perceived environment while *localizing* oneself in it. The two tasks are coupled in such a way so as to benefit each other; a good localization is crucial to create good maps, and a good map is necessary for localization. For this reason, the two tasks must be performed *simultaneously*, and hence, the full acronym SLAM. In recent years, the maturity of both online SLAM algorithms, together with fast and reliable image processing tools from the computer vision literature, has crystallized into a considerable quantity of real-time demonstrations of visual SLAM.

In this paper, we insist on the quality of the achieved localization, which will impact in turn the map quality. The key to good localization is to ensure the correct processing of the geometrical information gathered by the cameras. In this long introduction, we present an overview of visual SLAM and related techniques to show that visual SLAM systems have historically discarded

Manuscript received June 15, 2007; revised May 8, 2008. First published xxx; current version published xxx. This paper was recommended for publication by Associate Editor J. Tardos (with approval of the Guest Editors) and Editor L. Parker upon evaluation of the reviewers' comments.

The authors are with the Laboratoire d'Analyse et d'Architecture des Systèmes, Centre National de la Recherche Scientifique (LAAS-CNRS), University of Toulouse, Toulouse 31077, France (e-mail: jsola@laas.fr; monin@laas.fr; michel@laas.fr; tvidal@laas.fr).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TRO.2008.2004640

precious sensory information. We present a novel approach that uses the SLAM filter as a classical fusion engine that incorporates the full monocular information coming from multiple cameras.

A. Monocular SLAM

Possibly, the best example of the aforementioned technological crystallization is monocular SLAM, a particular case of bearings-only (BO) SLAM (where the sensor does not provide any range or depth). It is well known that the reduction in system observability due to BO measurements has two main drawbacks: the loss of the scale factor and the delay in obtaining good 3-D estimates. Previous works either added some metric measurement to observe the scale factor, such as odometry [1] or the size of known perceived objects [2], [3], or have considered it irrelevant [4]. The delay in getting good 3-D estimates comes from the fact that such estimates require several BO observations from different viewpoints. This makes landmark initialization in BO-SLAM difficult, to the point that satisfactory methods able to exploit all the geometrical information provided by the cameras have only recently become available. We have witnessed an evolution of the algorithms as follows. First, *delayed landmark initialization* methods attempted to obtain a full 3-D estimate before initialization via several observations from different viewpoints. Davison [3] showed real-time feasibility of monocular SLAM with affordable hardware, using the original extended Kalman filter (EKF) SLAM algorithm for all but the unmeasured landmark's depth, and a separate particle filter to estimate this depth. Initialization was *deferred* to the moment when the depth estimate was good enough. The consequence of a delayed scheme is that we can only initialize landmarks with enough parallax, i.e., those that are close to the camera and situated perpendicularly to its trajectory, and therefore, the need to operate in room-size scenarios with lateral motions. Second, Solà *et al.* [1] showed that *undelayed landmark initialization* (mapping the landmarks from their first, partial observation) was needed when considering low parallax landmarks, i.e., those that are remote and/or situated close to the motion axis. This permits mapping larger scenes while performing frontal trajectories. Third, Civera *et al.* [5] have recently achieved the mapping of landmarks up to infinity, due to an undelayed initialization via an *inverse depth parameterization* (IDP). IDP has also been developed by Eade *et al.* [6] in a FastSLAM2.0 context. Today, the monocular SLAM systems exploit the geometrical information in its entirety: from the first observation, independently of the sensor's trajectory, and up to the infinity range.

90 B. Structure From Motion (SFM)

91 Monocular SLAM compares to a similar problem solved
92 by the vision community: the structure from motion problem
93 (SFM). In SFM, the goal is to determine, from a collection of
94 images and up to an unrecoverable scale factor, the 3-D structure
95 of the perceived scene and all 6-D camera poses from where the
96 images were captured. When compared to SLAM, the structure
97 plays the role of the map, while the set of camera poses defines
98 all the successive observer's localizations.

99 Roboticians often claim that the main difference between
100 SFM and SLAM is that the former is solved offline via
101 the iterative nonlinear optimization method known as bundle
102 adjustment (BA) [7], while the latter must be incremen-
103 tally solved online, thus making use of stochastic estimators
104 or *filters* that naturally provide incremental operation. This
105 has been true for some years (today, SLAM is also solved
106 online with iterative optimization [8]), but does not tell the
107 whole story. The differences between SFM and SLAM are
108 not only in the methods but also in the objectives, meaning
109 that similar aspects of similar problems are given different
110 priorities.

111 In particular, SFM exploits the visual information in its en-
112 tirety without the difficulties encountered in monocular SLAM.
113 Let us try to understand this curious fact. SFM puts the struc-
114 ture as a final objective, i.e., as a result of the whole process,
115 and the emphasis is placed on minimizing the errors in the
116 *measurement space*, thus using all the measured information.
117 On the other hand, the SLAM map has a central role, with
118 some of the operations (and particularly landmark initializa-
119 tion) being performed in map space, which is the system's *state*
120 *space*. The fact that this state space is not statically observable,
121 because it is of higher dimension than the observation space,
122 leads to the difficulties exposed before. As an informal attempt
123 to fill this gap, we could say that modern undelayed methods
124 for monocular SLAM, with partial landmark initialization and
125 partial updates, are almost equivalent to an operation in the
126 measurement space: the information is initialized in the map
127 space *partially*, i.e., exactly as it comes from the measurement
128 space. A similar point of view over this concept can be found
129 in [9].

130 C. Stereovision SLAM

131 Stereovision SLAM has also received considerable attention.
132 The ability of a stereo assembly to directly and immediately pro-
133 vide 3-D landmark estimates allows us to use the best available
134 SLAM algorithms and rapidly obtain good results with little
135 effort in the conceptual parts. Such SLAM systems consider
136 the stereo assembly as being a single monolithic sensor, capa-
137 ble of gathering 3-D geometrical information from the robot's
138 surroundings, e.g. [10]. This fact, which appears perfectly rea-
139 sonable, is the main paradigm that this paper questions. By
140 considering two linked cameras as a single 3-D sensor, SLAM
141 is unable to face the following two issues.

142 1) *Limited 3-D Estimability Range*: While cameras are ca-
143 pable of sensing visible objects that are potentially at infinity,
144 a stereo rig provides only reasonably good 3-D estimates up

to a limited range, typically from 3 m to a few tens of meters 145
depending on the baseline. Because classical, nonmonocular 146
SLAM algorithms expect full 3-D estimates for landmark initial- 147
ization (i.e., they are reasoned in the map space), information 148
belonging to only this limited region can be used for SLAM. 149
This is really a pity; it is like if, having our two eyes, we were 150
obliged to neglect everything outside a certain range from us, 151
what we could call "*walking inside dense fog*." Without remote 152
landmarks, it is easy to lose spacial references, to become disori- 153
ented, and finally, find ourselves lost. Therefore, stereovision, 154
as it is classically conceived, is a bad starting point for visual 155
SLAM. 156

2) *Mechanical Fragility*: If we aim at extending the 3-D 157
estimability range beyond these few tens of meters, we need 158
to increase the stereo baseline while keeping or improving the 159
overall sensor precision. This is obviously a contradiction: larger 160
assemblies are less precise when using the same mechanical 161
solutions. In order to maintain accuracy with a larger assembly, 162
we must use more complex structures that will be either heavier 163
or more expensive, if not both. The result for moderately large 164
baselines (>1 m) is a sensor that is very easily decalibrated, 165
and therefore, almost useless. Large rigs, however, are very 166
interesting in outdoor applications because they allow farther 167
objects to be positioned, thus making them contribute to the 168
observability of the overall scale factor. This is especially true 169
in aerial and underwater settings where, without nearby objects 170
to observe, a small stereo rig provides no significant gain with 171
respect to a single camera. Self-calibration can compensate for 172
the inherent lack of stability of large camera rigs. It also allows 173
multicamera platforms to start operation without undergoing a 174
previous calibration phase, making on-field system deployment 175
and maintenance easier. 176

To our knowledge, the only SLAM work that goes beyond the 177
current stereoparadigm (apart from our conference paper [11]) 178
is the one by Paz *et al.* [12], which uses a small-baseline, fully 179
calibrated stereo rig. Matched features presenting significant 180
disparity are initialized as classical Euclidean landmarks, while 181
those presenting low disparities are treated with the inverse 182
depth algorithm. 183

184 D. Visual Odometry (VO)

185 One could say that, in terms of methodology, visual odom- 186
etry (VO) is to stereovision SLAM what SFM is to monocular 187
SLAM. VO is conceived to obtain the robot's ego motion from 188
a sequence of stereo images [13]. Visual features are matched 189
across two or more pairs of stereo images taken during the robot 190
motion. An iterative minimization algorithm, usually based on 191
BA, is run to recover the stereo rig motion, which is then trans- 192
formed into robot motion. For this, the algorithm needs to re- 193
cover the structure of the 3-D points that correspond to the 194
matched features. This structure is not exploited for other tasks 195
and can be usually discarded. Remarkably, when the structure 196
is coded in the measurement space (u, v, d) , a disparity $d \rightarrow 0$ 197
allows points at infinity to be properly handled [14]. This is also 198
accomplished by using homogeneous coordinates [7]. VO must 199
work in real time because robot localization is needed online.

200 Advanced VO solutions achieve very low drift levels after long
 201 distances by making use of: 1) hardware-based image process-
 202 ing with real-time construction and querying of large feature
 203 databases [15]; 2) dense image information matching via planar
 204 homographies and the use of the quadrifocal tensor [16]; or 3)
 205 bundle adjusting the set of N recent key frames together with
 206 additional fusion with an inertial measurement unit (IMU) [14].

207 E. Sensor Fusion in SLAM

208 The fact of SLAM being solved by filters allows us to envision
 209 SLAM systems as sensor fusion engines. Let us highlight some
 210 of the assets of filtering in sensor fusion.

- 211 1) *Multisensor operation*: Any number of differing sensors
 212 can be operated together in a consistent framework.
- 213 2) *Sensors self-calibration*: Unknown biases, gains, and
 214 other sensor's parameters can be estimated provided that
 215 they are observable [17].
- 216 3) *Desynchronized operation*: The data rates of all these sen-
 217 sors do not need to be synchronized.
- 218 4) *Decentralized operation*: Advanced filter formulations
 219 such as those using channel filters [18] achieve a decen-
 220 tralized operation that should permit live connection and
 221 disconnection of sensors without the need for filter repro-
 222 gramming or reparameterization.

223 This paper explores the first three points for the case of mul-
 224 tiple cameras.

225 SLAM systems naturally fuse information from both propri-
 226 oceptive (odometry, GPS, and IMU) and exteroceptive (range
 227 scanners, sonar, and vision) sensors into the map. But our in-
 228 terest here is in fusing several exteroceptive sensors. We can
 229 distinguish two cases.

- 230 1) *Sensors of different kind*: When using differing sensors
 231 (e.g., laser plus vision), the main problem is in finding a
 232 map representation well adapted to the different kinds of
 233 sensory data (i.e., the data association problem).
- 234 2) *Sensors of the same kind*: The perceived information is of
 235 the same nature. This makes appearance-based matching
 236 possible, and therefore, makes map building easier. Nev-
 237 ertheless, most of such SLAM systems do not take advan-
 238 tage of fusion. Instead, the extrinsic parameters linking
 239 the sensors are calibrated offline, and the set of sensors
 240 is treated as a single supersensor. This is the case for
 241 two 180° range scanners simulating a 360° one, and for
 242 the previously mentioned stereo rig simulating a 3-D sen-
 243 sor. A sensor-fusion approach in these cases should nat-
 244 urally bring the aforementioned advantages to the SLAM
 245 system.

246 F. Multicamera SLAM and the Aim of This Paper

247 The key idea of this paper is very simple: by employing
 248 the SLAM filter as a fusion engine, we will be able to use
 249 any number of cameras in any configuration. And, by treat-
 250 ing them as BO sensors with the modern undelayed initializa-
 251 tion methods, we will extract the entire geometrical information
 252 provided by the images. The filter—not the sensor—will be re-

sponsible for making the 3-D properties of the perceived world
 arise.

Applications may vary from the simplest stereo system,
 through robots with several differing cameras (e.g., a panoramic
 one for localization and a perspective one looking forward
 for reactive navigation), to multirobot cooperative SLAM
 where BO observations from different robots are used to
 determine the 3-D locations of very distant landmarks. Al-
 though there certainly exist issues concerning multicamera
 management, the main ideas we want to convey may be
 demonstrated with systems of just two cameras. In this pa-
 per, we will illustrate two cases: first, the case of a robot
 equipped with a stereo rig, with its cameras being treated
 as two individual monocular sensors and second, two cam-
 eras moving independently and mapping together an outdoors
 scene.

This paper draws on previous work published in the confer-
 ence paper [11] and the author's Ph.D. thesis [19]. These two
 works use the federated information sharing algorithm (FIS)
 in [1] to initialize the landmarks, which has been surpassed by
 the inverse depth methods (IDP) [5]. The present paper takes
 and extends all this research by developing a better founded jus-
 tification (providing a wider scope to the proposed concepts), by
 improving on the implementation with the incorporation of IDP
 in the algorithms, and by extending the experimental validation
 to a cooperative monocular SLAM setup.

This paper is organized as follows. Section II presents the
 main ideas that will be exploited later and revises some back-
 ground material for monocular SLAM. Section III explains how
 to set up multicamera SLAM, an application for stereo benches
 with self-calibration, and an application for two collaborative
 cameras. Section IV presents the perception and map manage-
 ment techniques used. Sections V and VI show the experimen-
 tal results, and finally, Section VII gives conclusions and future
 directions.

II. 3-D ESTIMABILITY IN VISUAL SLAM

In this section, we present the ideas that support our approach
 to visual SLAM. We make use of the concept of estimability,
 which will help understand the abilities of vision for observing
 3-D structure in the presence of uncertainty. We clarify the key
 properties of undelayed initialization in monocular SLAM, and
 remark its importance in multicamera SLAM. We also remind
 the key aspects of IDP-SLAM.

A. Geometrical Approach to 3-D Estimability

We are interested in finding the shape and dimensions of the
 3-D-estimable region defined by two monocular views.

For this, we start with a couple of ideas to help understand-
 ing the concept of estimability used. When a new feature is
 detected in an image, the backprojection of its noisy-measured
 position defines a conic-shaped *pdf* for the landmark position,
 called *ray*, which extends to infinity (see Fig. 1). Let us con-
 sider two features extracted and matched from a pair of images,
 corresponding to the same landmark: their backprojections are
 two conic rays A and B that extend to infinity. The angular

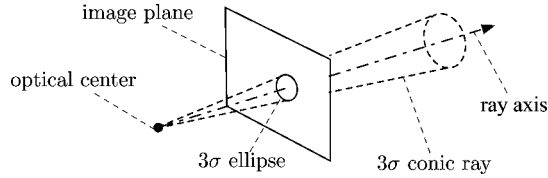


Fig. 1. Conic ray backprojects the elliptic representation of the Gaussian 2-D measure. It extends to infinity.

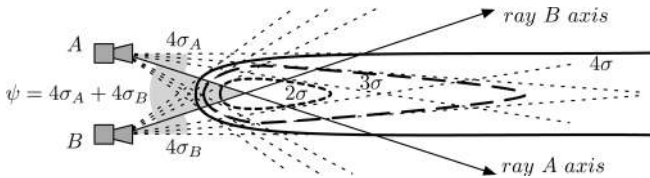


Fig. 2. Different regions of intersection for (solid) 4σ , (dashed) 3σ , and (dotted) 2σ ray widths when the outer 4σ bounds are parallel. (Shaded) The parallax or angle between rays axes A and B is $\psi = 4\sigma_A + 4\sigma_B$.

widths of these rays can be defined as a multiple of the standard deviations σ_A and σ_B of the angular errors (a composition of the cameras extrinsic and intrinsic parameters errors, and of the image processing algorithms accuracy). Informally speaking, we may say that the landmark's depth is fully estimated if the region of intersection of these rays is both *closed* and *sufficiently small*. If we consider, for example, the case where the two external 4σ bounds of the rays are parallel (see Fig. 2), then we can assure that the 3σ intersection region (which covers 98% probability) is *closed* and that the 2σ one (covering 74%) is *closed and small*. The ratio between the depth's standard deviation and its mean (a measure of linearity in monocular EKF-SLAM [1], [3]) is then better than 0.25. The *parallax* angle ψ between the two rays axes is therefore $\psi = 4(\sigma_A + \sigma_B) = \text{constant}$. This is the minimum parallax for full estimability.

In 2-D, we can plot the locus of constant estimability. In the case, where σ_A and σ_B can be considered constant, ψ is constant too, and from the inscribed angle theorem, the locus is then circular (Fig. 3, see also [19]). Landmarks inside this circle are considered *fully estimable*—and *partially* outside. In 3-D, the *fully 3-D estimable* region is obtained by revolution of this circle around the axis joining both cameras, producing a torus-shaped region with a degenerated central hole. This shape admits the following interpretations.

- 1) In a stereo configuration or for a lateral motion of a moving camera (see Fig. 3, left), the estimable region is located in front of the sensor. Beyond the region's border stereo provides no profit: if we want to consider distant landmarks, we have to use undelayed monocular techniques.
- 2) Depth recovery is impossible in the motion axis of a single camera moving forward (Fig. 3, right). Close to this axis, estimability is possible only if the region's radius becomes very large. This implies the necessity of very large displacements of the camera during the initializa-

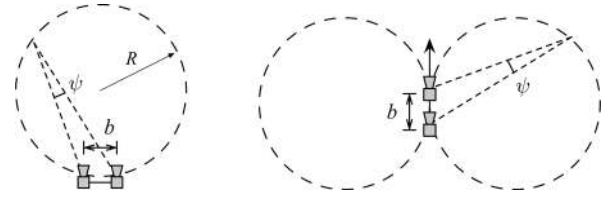


Fig. 3. Simplified depth estimability regions in a (left) stereo rig and (right) a camera traveling forward. The angle ψ is the one that assures estimability via triangulation from different viewpoints. The maximum range is $2R = b/\sin(\psi/2)$.

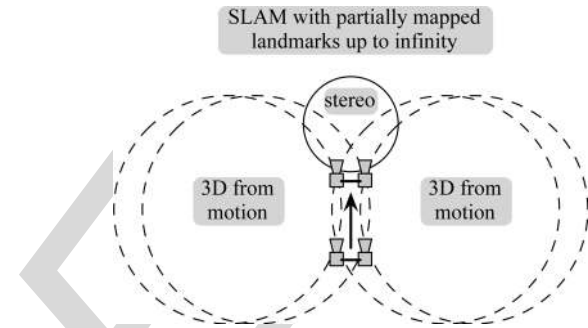


Fig. 4. Simplified depth estimability for a stereo rig moving forward. On both sides, estimability depends on the baseline gained by motion. In front, by stereo. Out of these bounds and up to infinity, landmarks are mapped partially. SLAM keeps incorporating the visual information due to the undelayed monocular methods, i.e., IDP in our case.

tion process. Again, this can be accomplished only with undelayed initializations.

- 3) By combining both monocular and stereovision, we get an instant estimability of close frontal objects while still utilizing the information of distant ones (see Fig. 4). Landmarks lying outside the estimability regions are not 3-D-estimable but, when initialized using undelayed monocular methods, they will contribute to constrain the camera orientation. Ideally, long-term observations of stable distant landmarks would completely cancel orientation drift (visual compass).

B. Monocular IDP-SLAM

The core algorithm of this paper is an EKF-SLAM with an IDP of landmarks during the initialization phase, as described in [5]. In IDP-SLAM, partially observed landmarks are coded as a 6-D-vector,

$$\mathbf{i} = [\mathbf{x}_0, \theta, \psi, \rho] \quad (1)$$

where \mathbf{x}_0 is the 3-D position of the camera at initialization time, (θ, ψ) are the elevation and azimuth angles in global frame defining the direction of the landmark's ray, and ρ is the inverse of the Euclidean distance from \mathbf{x}_0 to the landmark's position (notice that ρ is usually known as *inverse depth* but it is rather an inverse distance). After the first observation, all parameters of \mathbf{i} except ρ are immediately observable, and their values and covariances are obtained by proper inversion and linearization of the observation functions. The inverse depth ρ is initialized

369 with a Gaussian $\mathcal{N}(\rho - \bar{\rho}; \sigma_\rho^2)$ such that in the depth dimension
370 $s = 1/\rho$, we have

$$s_{(-n\sigma)} = \frac{1}{\bar{\rho} - n\sigma_\rho} = \infty \quad (2)$$

$$s_{(+n\sigma)} = \frac{1}{\bar{\rho} + n\sigma_\rho} = s_{\min} \quad (3)$$

371 with s_{\min} the minimum considered depth and n the inverse depth
372 shape factor. This gives $\bar{\rho} = 1/(2s_{\min})$ and, more remarkably

$$n\sigma_\rho = \bar{\rho}. \quad (4)$$

373 Importantly, values of $1 \leq n \leq 2$ assure from (2) that the infinity
374 range is included in the parametrization with ample probability.

375 On subsequent updates, IDP achieves correct EKF operation
376 (i.e., quasi-linear behavior) along the whole ray as long as the
377 parallax shown by the new viewpoint is not too large. The linear-
378 earity test in [20] is regularly evaluated. If passed, the landmark
379 can be safely transformed into a 3-D Euclidean parametrization.

380 III. MULTICAMERA SLAM

381 The general scheme for the multicamera SLAM system is
382 presented in this section. This scheme is particularized to deal
383 with two different problems. The first one is the automatic self-
384 calibration of a stereo rig while performing SLAM. The second
385 one is a master-lave solution to cooperative monocular SLAM.
386 Both setups are explained here, and their corresponding experi-
387 ments are presented in Sections V and VI.

388 A. System Overview

389 We implement the multicamera SLAM system as follows. A
390 central EKF-SLAM will hold the stochastic representation of
391 the set of all cameras \mathcal{C}_i plus the set of landmarks \mathcal{L}_j

$$X^\top = [\mathcal{C}_1^\top \quad \dots \quad \mathcal{C}_N^\top \quad \mathcal{L}_1^\top \quad \dots \quad \mathcal{L}_M^\top] \quad (5)$$

392 where the cameras states contain position and orientation quaternion
393 [$\mathcal{C}_i = (\mathbf{r}_i, \mathbf{q}_i) \in \mathbb{R}^7$], and landmarks can be coded either
394 in inverse depth ($\mathcal{L}_j = \mathbf{i}_j \in \mathbb{R}^6$) or in Euclidean coordinates
395 ($\mathcal{L}_j = \mathbf{p}_j \in \mathbb{R}^3$). Any number of cameras can be considered
396 this way. As each camera needs to remain localized properly,
397 it needs to observe a minimum number of landmarks at each
398 frame. The algorithm's complexity increases linearly with the
399 number of cameras if this number is small with respect to the
400 map.

401 For camera motions, we consider two possible models. In
402 the first one, a simple odometer provides motion predictions
403 $[\Delta x, \Delta y, \Delta \psi]$ in the robot's local 2-D plane. Gaussian uncer-
404 tainties are added to the 6-DOF linear and angular components
405 $[x, y, z, \phi, \theta, \psi]$ with a variance proportional to the measured
406 forward motion Δx

$$\{\sigma_x^2, \sigma_y^2, \sigma_z^2\} = k_L^2 \cdot \Delta x \quad (6)$$

$$\{\sigma_\phi^2, \sigma_\theta^2, \sigma_\psi^2\} = k_A^2 \cdot \Delta x. \quad (7)$$

407 The variance in $[\phi, \theta, \psi]$ is mapped to the quaternion space using
408 the corresponding Jacobians.

The second model is a 6-DOF constant velocity model

$$\mathbf{r}^+ = \mathbf{r} + \mathbf{v} \Delta t$$

$$\mathbf{q}^+ = \mathbf{q} \times v2q(\boldsymbol{\omega} \Delta t)$$

$$\mathbf{v}^+ = \mathbf{v} + \boldsymbol{\eta}_v$$

$$\boldsymbol{\omega}^+ = \boldsymbol{\omega} + \boldsymbol{\eta}_\omega$$

where $()^+$ means the updated value, \times is the quaternions prod- 410
uct, and $v2q(\boldsymbol{\omega} \Delta t)$ transforms the local incremental rotation 411
vector $\boldsymbol{\omega} \Delta t$ into a quaternion (quaternions are systematically 412
normalized). This way, the camera state vector \mathcal{C}_i is augmented 413
to $\mathcal{C}_i = (\mathbf{r}_i, \mathbf{q}_i, \mathbf{v}_i, \boldsymbol{\omega}_i) \in \mathbb{R}^{13}$. At each time step, perturbations 414
 $\{\boldsymbol{\eta}_v, \boldsymbol{\eta}_\omega\} \sim \mathcal{N}(0; \{\sigma_v^2, \sigma_\omega^2\})$ add variances to the linear and 415
angular velocities proportionally to the elapsed time Δt 416

$$\sigma_v^2 = k_v^2 \cdot \Delta t \quad (8)$$

$$\sigma_\omega^2 = k_\omega^2 \cdot \Delta t. \quad (9)$$

The events of camera motion, landmark initialization, and 417
landmark observation are handled as in regular IDP-SLAM by 418
just selecting the appropriate block elements from the SLAM 419
state vector and covariances matrix, and applying the corre- 420
sponding motion or observation models. For example, at the 421
observation of landmark j from camera i , we would use the 422
function $\mathbf{u}_j^i = \mathbf{h}(\mathcal{C}_i, \mathcal{L}_j)$, which will be explained later for the 423
case of an IDP ray [see 11]. Before transforming IDP rays into 424
points, the linearity test in [20] needs to hold for all cameras. 425

426 B. Stereo SLAM With Extrinsic Self-Calibration

Our approach is relevant to fully calibrated stereo rigs if they 427
are small (10–20 cm, as in [12]) or if, having long baselines, their 428
main extrinsic parameters can be continuously self-calibrated. 429

Not all of the six extrinsic parameters of a stereo rig (three for 430
translation, three for orientation) need to be calibrated. In fact, 431
the notion of *self-calibration* inherently requires the system to 432
possess its own gauge. In our case, the metric dimensions or 433
scale factor of the whole world-robot system can only be ob- 434
tained either from the stereo rig baseline, which is one of the 435
extrinsic parameters (then, it makes no sense to self-calibrate 436
the gauge), or from odometry, which is often much less accurate 437
than any coarse measurement we could make of this baseline. 438
Additionally, as cameras are actually angular sensors, vision 439
measurements are much more sensitive to the cameras orienta- 440
tions than to any translation parameter. This means that vision 441
measurements will contain little information about these trans- 442
lation parameters. In consequence, self-calibration may concern 443
only orientation, and more precisely, the orientation of one cam- 444
era with respect to the other. The error of the reconstructed map's 445
scale factor will be the same as the relative error of the baseline 446
measurement. 447

448 With these assumptions, our self-calibration solution is
449 straightforward: for the second camera, we just include its ori-
450 entation in the map and let EKF make the rest. The state vector
451 (5) is modified and written as

$$X^\top = [\mathcal{R}^\top \quad \mathbf{q}_R^\top \quad \mathcal{L}_1^\top \quad \dots \quad \mathcal{L}_M^\top]$$

452 where \mathcal{R} and $\mathcal{L}_1 \cdots \mathcal{L}_M$ are the robot pose and landmarks map.
 453 The left camera pose \mathcal{C}_L has a fixed transformation with respect
 454 to the robot, and \mathbf{q}_R is the orientation part of the right-hand
 455 camera \mathcal{C}_R in the robot frame. The time-evolution function of
 456 the angular extrinsic parameters is simply $\mathbf{q}_R^+ = \mathbf{q}_R + \gamma$, where
 457 γ is a white, Gaussian, low-energy process noise that accounts
 458 for eventual decalibrations, e.g., due to vibrations. For short-
 459 duration experiments, we set $\gamma = 0$. A coarse analysis of the
 460 stereo structure's mechanical precision will be enough to set the
 461 initial uncertainty to a value of the order of 1° or 2° per axis.
 462 This can be reduced to a few tenths of degree in cases where we
 463 dispose of previous calibrated values about which we are not
 464 confident anymore.

465 C. Cooperative Multicamera SLAM

466 The ideal, most general case of cooperative SLAM (5), corre-
 467 sponds to a (not too large) number of cameras moving indepen-
 468 dently. Each camera is able to manage its own measurements
 469 and communicates directly with the map. The aim of this com-
 470 munication is to obtain information about existing landmarks
 471 to get localized, and provide information about new or reob-
 472 served landmarks. This way, the algorithms to be executed by
 473 each camera are absolutely symmetrical, without any kind of
 474 hierarchy. A simplified implementation considers cameras with
 475 different privileges.

476 In our particular case, the cooperative SLAM system consid-
 477 ers two cameras. One of them takes the role of *master*, and
 478 is responsible for all landmarks detection and initialization.
 479 The second one acts as the *slave*. It follows the master at a
 480 close distance and reobserves the SLAM map that is being
 481 built by the master. By doing so, it provides a second view-
 482 point to landmarks just initialized, accelerating the convergence
 483 of the map. The master and slave trajectories are highly in-
 484 dependent, and for instance, they can cross paths. The only
 485 requirement is to look in the same direction. A trivial exten-
 486 sion to more than two cameras consists in including additional
 487 slaves.

488 IV. PERCEPTION AND MAP MANAGEMENT

489 Active search (AS, nicely described in [21] and also referred
 490 to as *top-down* in [6]) is a powerful framework for real-time
 491 image processing within SLAM. It has been successfully used in
 492 several monocular SLAM works [3], [5], [11], using a diversity
 493 of techniques for landmark initialization. The idea of AS is to
 494 exploit the information contained in the map to predict a number
 495 of characteristics of the landmarks to observe. AS is helpful in
 496 solving the following issues:

- 497 1) selecting interesting image regions for initialization;
- 498 2) selecting the most informative landmarks to measure;
- 499 3) predicting where in the image they may be found, and with
 500 which probability;
- 501 4) predicting the current landmark's appearance to maximize
 502 the chances of a successful match.

A. Feature Detection and Initialization

503

504 Based on the projection of the map information into the master
 505 image, a heuristic strategy is used to select a region of interest
 506 for a new initialization: we divide the image with a grid and
 507 randomly select a grid element with no landmarks inside. We
 508 extract the strongest Harris point [22] in this region and validate
 509 it if its strength is above a predefined threshold. We store a small
 510 rectangular region or *patch* of 15×15 pixels around the point
 511 as the landmark's appearance descriptor, together with the pose
 512 of the camera. Finally, we initialize the IDP ray in the SLAM
 513 map.

B. Expectations: The Active Search Regions

514

515 Some considerations about AS can be made for its usage in
 516 multicamera IDP-SLAM to improve performance. We use for
 517 this the \mathcal{E}_1 and \mathcal{E}_∞ ellipses, defined and explained as follows.

518 1) \mathcal{E}_1 *Ellipse: Expectation of the Inverse Depth Ray*: The
 519 inverse depth ray (1) is easily projected into a camera. We take
 520 the transformation to camera frame given in [5]:

$$\mathbf{h}_1^c = \mathbf{R}(\mathbf{q})^\top (\rho(\mathbf{x}_0 - \mathbf{r}) + \mathbf{m}(\theta, \psi)) \quad (10)$$

521 where $\mathbf{R}(\cdot)$ is the rotation matrix corresponding to the camera
 522 orientation \mathbf{q} and \mathbf{r} is the current camera position. This value
 523 is then projected into the camera, described by intrinsic and
 524 distortion parameters \mathbf{k} and \mathbf{d} (we use a classical radial distortion
 525 model of up to three parameters, which is inverted as explained
 526 in [19]). Let us call $\mathcal{K} = (\mathbf{k}, \mathbf{d})$ the camera parameters, $\mathcal{C} =$
 527 (\mathbf{r}, \mathbf{q}) the camera pose, and $\mathbf{i} = (\mathbf{x}_0, \theta, \psi, \rho)$ the IDP ray. The
 528 observation function is

$$\mathbf{u} = \mathbf{h}_1(\mathcal{C}, \mathcal{K}, \mathbf{i}) + \eta = \text{project}(\mathbf{h}_1^c, \mathcal{K}) + \eta \quad (11)$$

529 where $\text{project}(\cdot)$ takes into account the camera model (we use
 530 perspective cameras) and η is the pixel Gaussian noise, with
 531 covariance \mathbf{R} .

532 We define the \mathcal{E}_1 ellipse as the Gaussian expectation
 533 $\mathcal{E}_1(\mathbf{u}) \triangleq \mathcal{N}(\mathbf{u} - \bar{\mathbf{e}}_1; \mathbf{E}_1)$, with \mathbf{u} being the pixel position, and
 534 with mean and covariances matrix

$$\bar{\mathbf{e}}_1 = \mathbf{h}_1(\bar{\mathcal{C}}, \mathcal{K}, \bar{\mathbf{i}}) \quad (12)$$

$$\mathbf{E}_1 = [\mathbf{H}_c \mathbf{H}_i] \mathbf{P}_{c,i} [\mathbf{H}_c \mathbf{H}_i]^\top + \mathbf{R}. \quad (13)$$

535 Here, \mathbf{H}_c and \mathbf{H}_i are the Jacobians of \mathbf{h}_1 with respect to the
 536 uncertain parameters \mathcal{C} and \mathbf{i} , \bullet are variable estimates from
 537 the SLAM map, and $\mathbf{P}_{c,i}$ is the joint covariances matrix (all
 538 correlations and cross correlations) of \mathcal{C} and \mathbf{i} , also from the
 539 map. In AS, \mathcal{E}_1 is usually gated at 3σ , giving place to an elliptic
 540 region in the image where the landmark must project with 98%
 541 probability. However, this is not necessarily true in cases of
 542 noticeable parallax, as we examine now.

543 At landmark initialization, its inverse depth ρ is initialized
 544 according to (2)–(4). When considering 3σ uncertainty regions,
 545 (4) implies that ρ can go negative with a nonnegligible probab-
 546 ility, meaning that the coded landmarks might be situated *behind*
 547 *the camera*. This becomes evident when projecting the IDP ray
 548 into a second camera presenting some parallax: the projected
 549 3σ \mathcal{E}_1 ellipse contains a region with negative disparity (see

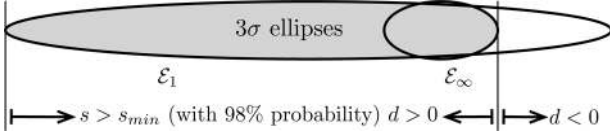


Fig. 5. 3σ search region defined by the \mathcal{E}_1 ellipse contains a significant part that corresponds to negative disparities $d < 0$, where the feature should not be searched. The final 3σ search region (gray) is defined by the \mathcal{E}_1 and \mathcal{E}_∞ ellipses. The rightmost 3σ border of \mathcal{E}_∞ is where the probability to find the projection of the infinity point has fallen below 2%.

Q1

550 Fig. 5). It is desirable to limit the search area to values of only
 551 positive disparity for two reasons: the correlation-based search
 552 (one of the most time-consuming processes) is faster and the
 553 possibility of including false matches as outliers is diminished.
 554 With nonrectified images and/or camera sets with uncertain ex-
 555 trinsic parameters, determining the null disparity bound is not
 556 straightforward. One solution is to use the \mathcal{E}_∞ ellipse, which we
 557 introduce in the following paragraph.

558 2) \mathcal{E}_∞ Ellipse: *Expectation of the Infinity Point*: The infinity
 559 point is easily projected by considering the transformation (10)
 560 with $\rho \rightarrow 0$

$$\mathbf{h}_\infty^c \approx \mathbf{R}(\mathbf{q})^\top \mathbf{m}(\theta, \psi) \quad (14)$$

561 where only the camera orientation \mathbf{q} and the ray's direction
 562 angles (θ, ψ) are present (the visual compass). Proceeding as
 563 before, we obtain the definition of the ellipse $\mathcal{E}_\infty(\mathbf{u}) \triangleq \mathcal{N}(\mathbf{u} -$
 564 $\bar{\mathbf{e}}_\infty; \mathbf{E}_\infty)$ as

$$\bar{\mathbf{e}}_\infty = \mathbf{h}(\bar{\mathbf{q}}, \bar{\mathcal{K}}, \bar{\theta}, \bar{\psi}) \quad (15)$$

$$\mathbf{E}_\infty = [\mathbf{H}_q \ \mathbf{H}_\theta \ \mathbf{H}_\psi] \mathbf{P}_{\{q, \theta, \psi\}} [\mathbf{H}_q \ \mathbf{H}_\theta \ \mathbf{H}_\psi]^\top + \mathbf{R} \quad (16)$$

565 where $\mathbf{P}_{\{q, \theta, \psi\}}$ is the joint covariances matrix of the uncertain
 566 parameters. The \mathcal{E}_∞ 3σ region is composed of the previous \mathcal{E}_1
 567 region, as indicated in Fig. 5, to define the search area.

568 C. Selection of the Best Map Updates

569 Following the AS approach in [23], a predefined number of
 570 landmarks with the biggest \mathcal{E}_1 ellipse surfaces are selected in
 571 each camera as those being the most interesting to be measured.
 572 For each camera, we organize all candidates (visible landmarks)
 573 in descending order of expectation surfaces, without caring if
 574 they are points or rays. We update at each frame a predefined
 575 number of them (usually around 10, and no more than 20).
 576 Updates are processed sequentially, with all Jacobians being
 577 recalculated each time to minimize the effect of linearization
 578 errors.

579 D. Feature Matching: Affine Patch Warping

580 AS continues by *warping* the stored patch and searching for
 581 a correlation peak inside the search area earlier. The objec-
 582 tive of warping is to predict the landmark's current appearance,
 583 maximizing the chances for a good match. In the absence of dis-
 584 tortion, a planar homography $\mathbf{H} \in \mathbb{R}^{3 \times 3}$, defined in the homo-
 585 geneous spaces, would be desirable [24]. This type of warping
 586 requires the online estimation of the patch normal in the 3-D

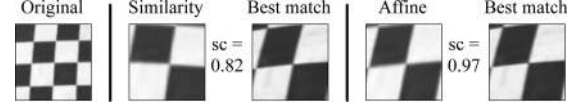


Fig. 6. Similarity and affine warping on a sample patch. From left to right: original patch; similarity warped patch ($\sim 180\%$ scale, 10° rotation); best match in a later image affected by distortion and its zero mean normalized cross correlation (ZNCC) score (0.82); affine warped patch; best match and score (0.97). The affine warping contains a significant skew component mainly due to image distortion. The improvement in the ZNCC score is very important.

587 space, and may become very time-consuming. A good simplifi-
 588 cation considers this normal fixed at the initial visual axis [23].
 589 Further simplification applies just a similarity transformation
 590 $\mathbf{T} = s\mathbf{R} \in \mathbb{R}^{2 \times 2}$ in the image Euclidean plane [19]. This ac-
 591 counts only for scale changes s and rotations \mathbf{R} obtained from
 592 the stored information (landmark position, camera initial, and
 593 current poses). However, in the presence of distortion, features
 594 lying close to the image borders suffer from additional deforma-
 595 tions. We developed a warping approach that easily adds a
 596 skew component to the operator \mathbf{T} (thus achieving fully affine
 597 warping, but not perspective warping; Fig. 6), based on the Ja-
 598 cobian of the function linking the first observation to the current
 599 one. Let us consider the backward observation model $\mathbf{g}(\cdot)$ for a
 600 camera A at initialization time $t = 0$, and the observation model
 601 $\mathbf{h}(\cdot)$ for a different camera B at current time $t \geq 0$

$$\mathbf{p} = \mathbf{g}(\mathcal{C}_A(0), \mathcal{K}_A, \mathbf{u}_A(0), s_A)$$

$$\mathbf{u}_B(t) = \mathbf{h}(\mathcal{C}_B(t), \mathcal{K}_B, \mathbf{p}).$$

602 Here, \mathbf{p} is the landmark's position, $\mathcal{K}_i = (\mathbf{k}_i, \mathbf{d}_i)$ are the intrin-
 603 sic and distortion parameters of camera i , $\mathbf{u}_i(t)$ is the measured
 604 pixel, and s_A is the landmark's depth with respect to the initial
 605 camera. We can compose these functions to obtain the expres-
 606 sion linking the initial and the current pixels
 607

$$\mathbf{u}_B(t) = \mathbf{h}[\mathcal{C}_B(t), \mathcal{K}_B, \mathbf{g}(\mathcal{C}_A(0), \mathcal{K}_A, \mathbf{u}_A(0), s_A)]. \quad (17)$$

608 When all but the pixel positions are fixed, this represents an
 609 invertible mapping $\mathbb{R}^2 \mapsto \mathbb{R}^2$ from the pixels in the first image
 610 to the pixels in the current one. The local linearization around
 611 the initially measured pixel defines an affine warping expressed
 612 by the Jacobian matrix

$$\mathbf{T} = \left. \frac{\partial \mathbf{u}_B}{\partial \mathbf{u}_A} \right|_{(\mathcal{C}_A(0), \mathcal{C}_B(t), \mathcal{K}_A, \mathcal{K}_B, \mathbf{u}_A(0), s_A)}. \quad (18)$$

613 By defining $\tilde{\mathbf{u}}_i$ as the coordinates of the patch in camera i , with
 614 the central pixel \mathbf{u}_i as the origin, we have $\tilde{\mathbf{u}}_B(t) = \mathbf{T} \tilde{\mathbf{u}}_A(0)$.
 615 Based on this mapping, we use linear interpolation of the pixels'
 616 luminosity to construct the warped patch.

617 V. EXPERIMENT 1: STEREO SLAM WITH SELF-CALIBRATION

618 The ‘‘White-board’’ indoor experiment aims at demonstrating
 619 stereovision SLAM with self-calibration. A robot with a stereo
 620 head looking forward is run for about 10 m in straight line inside
 621 the robotics laboratory at the LAAS (see Fig. 7). Over 500 image
 622 pairs are taken at approximately 5-Hz frequency. The robot

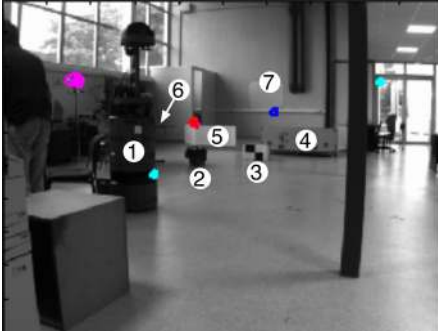


Fig. 7. Laboratoire d’Analyse et d’Architecture des System (LAAS) robotics laboratory. The robot will approach the scene in a straightforward trajectory. We notice in the scene the presence of a robot ①, a bin ②, a box ③, a trunk ④, a fence ⑤, a table ⑥ (hidden by the robot in this image), and the white board ⑦ at the end wall.

TABLE I
STEREO RIG PARAMETERS IN THE “WHITE-BOARD” EXPERIMENT

Scope	Parameters = Values
Dimensions	Baseline = 33 cm
Orientation - Euler	$\{\phi, \theta, \psi\} = \{0^\circ, 5^\circ, 0^\circ\}$
Cameras	$\{\text{resolution}, \text{FOV}\} = \{512 \times 384 \text{ pix}, 55^\circ\}$
Right camera uncertainties	$\{\sigma_\phi, \sigma_\theta, \sigma_\psi\} = \{1^\circ, 1^\circ, 1^\circ\}$

623 moves towards the objects to be mapped at 0.15 m/s. The stereo
624 rig consists of two intrinsically calibrated cameras arranged
625 as indicated in Table I. The orientations of both cameras are
626 specified with respect to the robot frame. The left camera is taken
627 as reference, thus deterministically specified, and the orientation
628 of the right one is initialized with an uncertainty of 1° standard
629 deviation. We use the odometry model (Section III-A) with
630 $k_L = 0.1 \text{ m}/\sqrt{\text{m}}$ and $k_A = 0.05 \text{ rad}/\sqrt{\text{m}}$.

631 We show details and results on the self-calibration procedure
632 and the metric accuracy of the resulting map. The mapping
633 process can be appreciated in the movie `whiteboard.mov` in
634 the multimedia section.

635 A. Self-Calibration

636 We plot in Fig. 8 left the evolution of the three self-calibrated
637 angles. We have also used the shape of the \mathcal{E}_∞ ellipses to pro-
638 vide additional qualitative evidence of the calibration process
639 (Fig. 9 and movie `whiteboard - e1nf.mov`). We observe the
640 following behavior.

641 1) *Pitch* θ : The pitch angle (cameras tilt, 5° nominal value) is
642 observable from the first matched landmark. It rapidly converges
643 to an angle of 4.77° and remains very stable during the whole
644 experiment.

645 2) *Roll* ϕ : Roll angle is observable after at least two land-
646 marks are observed from the right camera. Once this condition
647 holds, convergence occurs relatively fast.

648 3) *Yaw* ψ : Yaw angle is very weakly observable because
649 it is coupled with the landmarks depths: both yaw angle and
650 landmark depth variations produce a similar uncertainty growth
651 in the right image. For this reason, yaw converges slowly, only
652 showing reasonable convergence after some 50 frames.

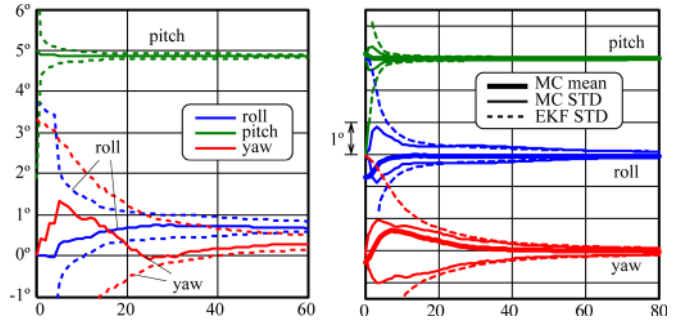


Fig. 8. Extrinsic self-calibration. (Left) The three Euler angles of the right camera orientation with respect to the robot during the first 60 frames. The 3σ bounds are plotted in dotted line showing consistent estimation. (Right) Error analysis after 100 MC runs using 200 frames per run (only the first 80 frames are shown). The thick solid lines represent the means over the 100 runs. The 3σ bounds for each angle are plotted using thin solid lines. The dotted lines represent the averaged 3σ bounds estimated by the EKF, showing consistent calibration.

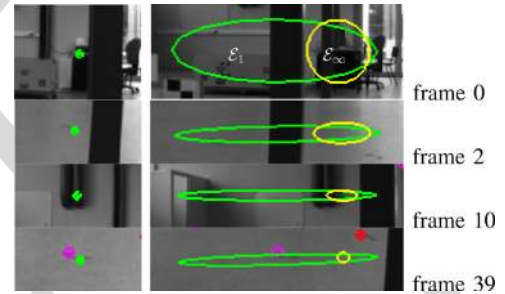


Fig. 9. Evolution of the \mathcal{E}_1 and \mathcal{E}_∞ ellipses during calibration. On the left column, newly detected pixels in the left image. On the right, expectations in the right image (green) \mathcal{E}_1 and (yellow) \mathcal{E}_∞ of the newly initialized IDP rays (i.e., still with the full initial uncertainty in ρ). At frame 0, initial uncertainties of 1° result in a big, round \mathcal{E}_∞ ellipse. After the first updated landmark from the left camera (frame 2), the uncertainty in pitch decreases and \mathcal{E}_∞ becomes flat. Successive updates further refine the calibrated angles. The yaw angle takes longer to converge, but the tiny \mathcal{E}_∞ in frame 39 shows that the calibration is already finished. The portion of the green ellipse on the right side of the yellow one corresponds to negative disparities and is not searched for matches. This portion is larger as parallax increases.

TABLE II
MC ANALYSIS OF THE SELF-CALIBRATION

Angle	MC mean	STD	EKF STD	Offline	STD
roll	0.69°	0.028°	0.018°	0.61°	0.013°
pitch	4.77°	0.003°	0.005°	4.74°	0.099°
yaw	0.33°	0.021°	0.016°	0.51°	0.109°

In Fig. 8 right, we plot results of a Monte Carlo (MC) analysis, run over the data of this experiment, for the mean and standard deviation of the Euler angles of the right camera. Because all MC runs are extracted from the same sequence, we tried to maximize their independence by using a different *random seed* in the algorithm (acting in the random selection of the initialization region, Section IV-A), and by starting each run at a different frame. The figure shows that the dynamic estimation is consistent (the EKF estimated sigmas are larger than the MC ones). After 200 frames, we compare these values with those of the offline calibration [25]. Table II summarizes these results, showing MC [(means and standard deviations (STD))] and Kalman Filter (EKF, showing the estimated STD). All

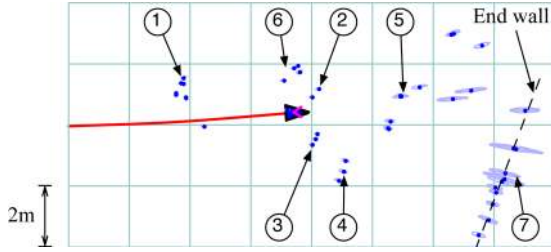


Fig. 10. Map produced during the “white board” experiment. We marked the mapped robot ①, the bin ②, the box ③, the trunk ④, the fence ⑤, the table ⑥, and the white board ⑦ at the end wall.

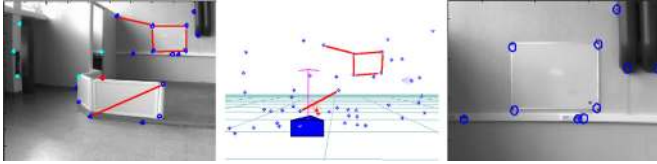


Fig. 11. Metric mapping. The magnitudes of some segments in the real laboratory are compared to those in the map (red lines). Ground truth corresponds to metric measurements of the distances between landmarks that are identified by zooming in the last image of the experiment (right) and translated to the real world. Thirteen points on the end wall are tested for coplanarity.

TABLE III
WHITE BOARD: MAP TO GROUND TRUTH TOMPARISON

segment	board	board	board	board	wall	fence
real (cm)	116	86	117	88	136	124
mapped	116.6	87.2	115.8	87.0	135.1	125.5
STD	0.91	0.81	1.21	0.52	1.06	1.32

666 self-calibrated values lie within the 3σ bounds defined by the
667 offline mean and STD values.

668 B. Metric Accuracy

669 We show in Fig. 10 a top view of the map generated during
670 this experiment. To contrast this map against reality, two tests
671 are performed: planarity and metric scale (see Fig. 11): 1) the
672 four corners of the white board are taken together with nine
673 other points at the end wall to test coplanarity: the 13 mapped
674 points are found to be coplanar within 4.9 cm STD; 2) the
675 lengths of the real and mapped segments marked in Fig. 11
676 are summarized in Table III. The white board has a physical
677 size of 120 cm \times 90 cm, but we take real measurements from
678 the approximated corners where the features are detected. We
679 observe errors in the order of 1 cm for landmarks that are still
680 about 4 m away from the robot.

681 VI. EXPERIMENT 2: COOPERATIVE MONOCULAR SLAM

682 This experiment shows independent cameras collaborating to
683 build a 3-D map using exclusively bearings-only observations.
684 Two independent cameras are placed on top of two bicycles
685 looking forward, moving on different trajectories in the park-
686 ing of the LAAS (see Fig. 12). Over 1000 images are taken
687 by each camera at 15-Hz frequency, 512×384 pixel resolution,
688 100° field of view (FOV), and are processed offline. The cam-



Fig. 12. Snapshots of master and slave sequences in cooperative SLAM. Faraway landmarks (e.g., black arrowed), still initialized as rays (red), are the ones fixing the orientation. Nearby landmarks, usually as Euclidean points (blue), maintain the metric. A virtual model of the master camera is visible from the slave camera (white arrowed). See *cooperativeSLAM.mov*.

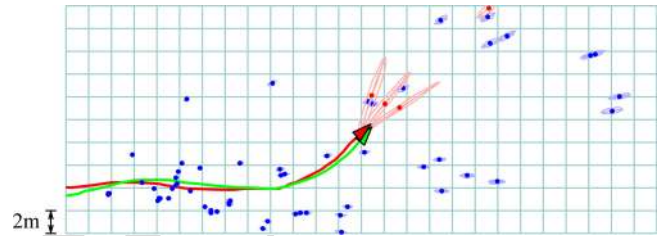


Fig. 13. Top view of the map produced by cooperative SLAM of two independent cameras, and their crossing trajectories. The grid spacing is 2 m.

eras travel approximately 28 m observing landmarks beyond 698
699 60 m. As in the previous experiment, the left camera is the master.
700 The two trajectories start parallel to each other, separated
701 75 cm perpendicularly to the motion direction. The reference
702 frame is defined by the master camera initial position and orientation,
703 which are initialized with null uncertainty. The scale factor is determined
704 by the initial baseline of 75 cm, meaning that the position of the slave
705 camera in the lateral Y -axis is also initialized with null uncertainty. The
706 orientations of the slave camera start with an uncertainty of 2° STD, and
707 its position in the frontal Y - and vertical Z -axes with $75 \text{ cm} \cdot \sin(2^\circ) = 2.6 \text{ cm}$
708 STD. With these uncertainties, the experiment’s initial configuration can
709 be set up manually by just observing the images and centering the projections
710 of some distant object. We use two independent constant-velocity models
711 with $k_v = 0.3 \text{ m/s} \cdot \sqrt{s}$ and $k_w = 0.3 \text{ rad/s} \cdot \sqrt{s}$. The measurement noise
712 is 1 pixel.

Landmarks at infinity, illumination changes and few salient 705
706 features are some characteristics of this outdoors scene. It presents
707 relatively few stable landmarks, something that makes the correct operation
708 of the SLAM system difficult. In the case of having crossing trajectories,
709 the problem of one camera occluding the other could appear and severely
710 affect the image processing. To avoid this, we decided to take both image
711 sequences shifted in time, i.e., one after the other, and make them overlap
712 for processing. The mapping process is presented in the movie *cooperativeSLAM.mov*
713 in the multimedia section. Fig. 13 shows the top view of the map and the
714 camera trajectories generated during this experiment.

A proper metrical evaluation of this experiment is difficult; 717
718 having a variable baseline does not allow us to compare the results, because
719 there is no knowledge of the ground truth. In order to evaluate this approach,
720 we consider the setup in experiment

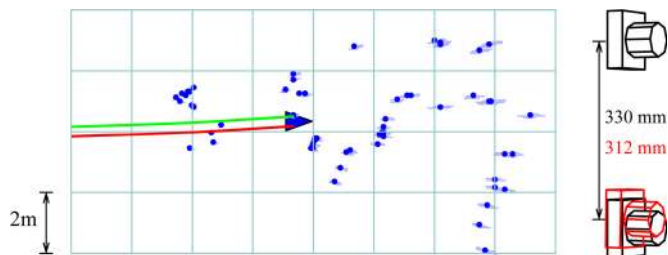


Fig. 14. Final map in the “white board” setup using the cooperative monocular SLAM algorithm. The cameras are modeled as being entirely independent using the same data and initial configuration as in Experiment 1. The stereo rig on the right shows (red) the final estimated relative position compared with (black) ground truth.

1 and apply the same algorithm. The new experiment consists of recovering the full extrinsic calibration, which is fixed in reality, considering both cameras as independent. Again, we use a constant-velocity model for each camera. The initial setup including uncertainties is as in experiment 1.

Fig. 14 shows the obtained map. We see that it compares very well to the map obtained in experiment 1 (see Fig. 10), where the motions of the two cameras were constrained by the stereo rig and a common motion was predicted using odometry. Fig. 14 bottom shows a detail of the cameras in their final relative position. We measure an error along the baseline of less than 2 cm. The orientation errors are less than 0.7° .

VII. CONCLUSION

We showed in this paper that fusing the visual information with monocular methods while performing multicamera SLAM provides several advantages: the ability to consider points at infinity, desynchronization of the different cameras, the use of any number of cameras of different types, sensor self-calibration, and the possibility to conceive decentralized schemes that will make realistic multirobot monocular SLAM possible. Except for decentralization, these advantages have been explored with the inverse depth monocular SLAM algorithm, and applied to two different problems: stereovision SLAM with an extrinsically decalibrated stereo rig and cooperative SLAM of two independently moving cameras.

Both demonstrations employed a *master-slave* approach, which made solving some of the issues of map and image management easier, and we are now improving on this by implementing a fully symmetrical approach. This approach should easily permit the extension of the presented applications to cases with more than two cameras. In parallel to these activities, we started new work on landmark parametrization to improve EKF linearity in cases of increasing parallax. Also, as parallax increases, landmarks appearances may change too much as to guarantee a stable operation with the matching methods presented here. We believe that wide baseline feature matching will be the bottleneck of visual SLAM for some time to come. As for decentralization, we note that it demands a full reformulation of the fusion engines we use in this paper (one central EKF), for example, via channel filters, and is currently a subject of intense research at LAAS and other laboratories.

REFERENCES

- [1] J. Solà, A. Monin, M. Devy, and T. Lemaire, “Undelayed initialization in bearing only SLAM,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Edmonton, AB, Canada, Aug. 2–6, 2005, pp. 2499–2504.
- [2] A. Chiuso, P. Favaro, H. Jin, and S. Soatto, “Structure from motion causally integrated over time,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 4, pp. 523–535, Apr. 2002.
- [3] A. J. Davison, “Real-time simultaneous localisation and mapping with a single camera,” in *Proc. Int. Conf. Comput. Vis.*, Nice, France, Oct. 2003, vol. 2, pp. 1403–1410.
- [4] J. Civera, A. J. Davison, and J. M. M. Montiel, “Dimensionless monocular SLAM,” in *Proc. Iberian Conf. Pattern Recognit. Image Anal.*, Jun. 2007, pp. 412–419.
- [5] J. Civera, A. Davison, and J. Montiel, “Inverse depth parametrization for monocular SLAM,” *IEEE Trans. Robot.*, vol. 24, no. 5, Oct. 2008.
- [6] E. Eade and T. Drummond, “Scalable monocular SLAM,” in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit.*, Jun. 17–22, 2006, vol. 1, pp. 469–476.
- [7] B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon, “Bundle adjustment—A modern synthesis,” in *Vision Algorithms: Theory and Practice*, ser. LNCS, W. Triggs, A. Zisserman, and R. Szeliski, Eds. New York: Springer-Verlag, 2000, pp. 298–375.
- [8] K. Konolige, “SLAM via variable reduction from constraints maps,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Barcelona, Spain, Apr. 18–22, 2005, pp. 667–672.
- [9] J. Folkesson, P. Jensfelt, and H. I. Christensen, “Vision SLAM in the measurement subspace,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Barcelona, Spain, Apr. 18–22, 2005, pp. 30–35.
- [10] J. Diebel, K. Reuterswård, S. Thrun, and R. G. J. Davis, “Simultaneous localization and mapping with active stereo vision,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Sendai, Japan, Oct. 2004, vol. 4, pp. 3436–3443.
- [11] J. Solà, A. Monin, and M. Devy, “BiCamSLAM: Two times mono is more than stereo,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Rome, Italy, Apr. 2007, pp. 4795–4800.
- [12] L. M. Paz, P. Piniés, J. Tardós, and J. Neira, “Large scale 6 DOF SLAM with stereo-in-hand,” *IEEE Trans. Robot.*, vol. 24, no. 5, Oct. 2008.
- [13] A. Mallet, S. Lacroix, and L. Gallo, “Position estimation in outdoor environments using pixel tracking and stereovision,” in *Proc. Int. Conf. Robot. Autom.*, San Francisco, CA, 2000, vol. 4, pp. 3519–3524.
- [14] K. Konolige, M. Agrawal, and J. Solà, “Large-scale visual odometry for rough terrain,” presented at the Int. Symp. Res. Robot., Hiroshima, Japan, Nov. 2007.
- [15] T. D. Barfoot, “Online visual motion estimation using FastSLAM with SIFT features,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Aug. 2–6, 2005, pp. 579–585.
- [16] A. I. Comport, E. Malis, and P. Rives, “Accurate quadrifocal tracking for robust 3D visual odometry,” in *Proc. Int. Conf. Robot. Autom.*, Rome, Italy, Apr. 10–14, 2007, pp. 40–45.
- [17] E. M. Foxlin, “Generalized architecture for simultaneous localization, auto-calibration, and map-building,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Lausanne, Switzerland, 2002, vol. 1, pp. 527–533.
- [18] E. Nettleton, H. Durrant-Whyte, and S. Sukkarieh, “A robust architecture for decentralised data fusion,” presented at the Int. Conf. Adv. Robot., Coimbra, Portugal, 2003.
- [19] J. Solà, “Towards visual localization, mapping and moving objects tracking by a mobile robot: a geometric and probabilistic approach” Ph.D. dissertation, Inst. Nat. Polytech. de Toulouse, Toulouse, France, 2007.
- [20] J. Civera, A. Davison, and J. Montiel, “Inverse depth to depth conversion for monocular SLAM,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Rome, Italy, Apr. 10–14, 2007, pp. 2778–2783.
- [21] A. J. Davison, “Active search for real-time vision,” in *Proc. Int. Conf. Comput. Vis.*, 2005, vol. 1, pp. 66–73.
- [22] C. Harris and M. Stephens, “A combined corner and edge detector,” in *Proc. 4th Alvey Vis. Conf.*, Manchester, U.K., 1988, pp. 189–192.
- [23] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, “MonoSLAM: Real-time single camera SLAM,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 6, pp. 1052–1067, Jun. 2007.
- [24] N. Molton, A. J. Davison, and I. Reid, “Locally planar patch features for real-time structure from motion,” presented at the Brit. Mach. Vis. Conf., Kingston, U.K., 2004.
- [25] K. Strobl, W. Sepp, S. Fuchs, C. Paredes, and K. Arbter. (2006). Camera calibration toolbox for Matlab. Inst. Robot. Mechatronics, Wessling, Germany, Tech. Rep. [Online]. Available: http://www.vision.caltech.edu/bouguetj/calib_doc/index.html

836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851

Joan Solà was born in Barcelona, Spain, in 1969. He received the B.Sc. degree in telecommunications and electronic engineering from the Universitat Politècnica de Catalunya, Barcelona, in 1995, the M.Sc. degree in control systems from the École Doctorale Systèmes, Toulouse, France, in 2003, and the Ph.D. degree in control systems from the Institut National Polytechnique de Toulouse in 2007, where he was hosted by the Laboratoire d'Analyse et d'Architecture des System (LAAS), Centre National de la Recherche Scientifique (CNRS).

He was a Postdoctoral Fellow at SRI International, Menlo Park, CA. He is currently at LAAS-CNRS, where he is engaged in research on visual localization and mapping. His current research interests include estimation and data fusion applied to off-road navigation, mainly using vision.

852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867

André Monin was born in Le Creusot, France, in 1958. He received the Graduate degree from the Ecole Nationale Supérieure d'Ingénieurs Electriciens de Grenoble, Grenoble, France, in 1980, the Ph.D. degree in nonlinear systems representation from the University Paul Sabatier, Toulouse, France, in 1987, and the Habilitation pour Diriger des Recherches degree from the University Paul Sabatier in 2002.

From 1981 to 1983, he was a Teaching Assistant with the Ecole Normale Supérieure de Marrakech, Marrakech, Morocco. Since 1985, he has been with the Laboratoire Automatique et d'Analyse des Systèmes, Centre Nationale de la Recherche Scientifique (LAAS-CNRS), Toulouse, as the "Chargé de Recherche." His current research interests include the areas of nonlinear filtering, systems realization, and identification.



Michel Devy received the degree in computer science engineering from the Ecole Nationale Supérieure d'Ingénieurs Electriciens de Grenoble, Grenoble, France, in 1976 and the Ph.D. degree from the Laboratoire Automatique et d'Analyse des Systèmes, Centre Nationale de la Recherche Scientifique (LAAS-CNRS), Toulouse, France, in 1980.

Since 1980, he has been with the Department of Robotics and Artificial Intelligence, LAAS-CNRS, where he is the Research Director and the Head of the Research Group Robotics, Action, and Perception. His current research interests include computer vision for automation and robotics applications. He has also been involved in numerous national and international projects concerning, about field and service robots, 3-D vision for intelligent vehicles, 3-D metrology, and others. He has authored or coauthored about 150 scientific communications.

868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884

Q3

Q2



Teresa Vidal-Calleja received the B.Sc. degree in mechanical engineering from the Universidad Nacional Autónoma de México, México City, México, in 2000, the M.Sc. degree in mechatronics from CINVESTAV-IPN, México City, in 2002, and the Ph.D. degree in robotics, automatic control, and computer vision from the Universitat Politècnica de Catalunya, Barcelona, Spain, in 2007.

She was a Visiting Research Student with the University of Oxford's Robotics Research Group, the Australian Centre for Field Robotics, and the University of Sydney. She is currently a Postdoctoral Fellow with the Robotics and Artificial Intelligence Group, Laboratoire Automatique et d'Analyse des Systèmes, Centre Nationale de la Recherche Scientifique (LAAS-CNRS), Toulouse, France. Her current research interests include autonomous vehicles, perception, control, and cooperation.

885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901

QUERIES

- 903 Q1: Author: Please reframe the references to colors in the caption of Figs. 5, 8, 9, 11, 12, and 14, if the artwork is not being
904 produced in color
- 905 Q2. Author: Please check if the details of the academic degrees received by A. Monin are OK as edited.
- 906 Q3. Author: Please provide the title of first degree. Also, provide the subject (physics, mathematics, electrical engineering, etc.)
907 in which M. Dey received the Ph. D. degree.

IEEE
PROOF

Fusing Monocular Information in Multicamera SLAM

Joan Solà, André Monin, Michel Devy, and Teresa Vidal-Calleja

Abstract—This paper explores the possibilities of using monocular simultaneous localization and mapping (SLAM) algorithms in systems with more than one camera. The idea is to combine in a single system the advantages of both monocular vision (bearings-only, infinite range observations but no 3-D instantaneous information) and stereovision (3-D information up to a limited range). Such a system should be able to instantaneously map nearby objects while still considering the bearing information provided by the observation of remote ones. We do this by considering each camera as an independent sensor rather than the entire set as a monolithic supersensor. The visual data are treated by monocular methods and fused by the SLAM filter. Several advantages naturally arise as interesting possibilities, such as the desynchronization of the firing of the sensors, the use of several unequal cameras, self-calibration, and cooperative SLAM with several independently moving cameras. We validate the approach with two different applications: a stereovision SLAM system with automatic self-calibration of the rig's main extrinsic parameters and a cooperative SLAM system with two independent free-moving cameras in an outdoor setting.

Index Terms—Calibration, image sequence analysis, Kalman filtering, machine vision, robot vision systems, stereovision.

I. INTRODUCTION

THE SIMULTANEOUS localization and mapping (SLAM) problem, as formulated by the robotics community, is that of creating a *map* of the perceived environment while *localizing* oneself in it. The two tasks are coupled in such a way so as to benefit each other; a good localization is crucial to create good maps, and a good map is necessary for localization. For this reason, the two tasks must be performed *simultaneously*, and hence, the full acronym SLAM. In recent years, the maturity of both online SLAM algorithms, together with fast and reliable image processing tools from the computer vision literature, has crystallized into a considerable quantity of real-time demonstrations of visual SLAM.

In this paper, we insist on the quality of the achieved localization, which will impact in turn the map quality. The key to good localization is to ensure the correct processing of the geometrical information gathered by the cameras. In this long introduction, we present an overview of visual SLAM and related techniques to show that visual SLAM systems have historically discarded

Manuscript received June 15, 2007; revised May 8, 2008. First published xxx; current version published xxx. This paper was recommended for publication by Associate Editor J. Tardos (with approval of the Guest Editors) and Editor L. Parker upon evaluation of the reviewers' comments.

The authors are with the Laboratoire d'Analyse et d'Architecture des Systèmes, Centre National de la Recherche Scientifique (LAAS-CNRS), University of Toulouse, Toulouse 31077, France (e-mail: jsola@laas.fr; monin@laas.fr; michel@laas.fr; tvidal@laas.fr).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TRO.2008.2004640

precious sensory information. We present a novel approach that uses the SLAM filter as a classical fusion engine that incorporates the full monocular information coming from multiple cameras.

A. Monocular SLAM

Possibly, the best example of the aforementioned technological crystallization is monocular SLAM, a particular case of bearings-only (BO) SLAM (where the sensor does not provide any range or depth). It is well known that the reduction in system observability due to BO measurements has two main drawbacks: the loss of the scale factor and the delay in obtaining good 3-D estimates. Previous works either added some metric measurement to observe the scale factor, such as odometry [1] or the size of known perceived objects [2], [3], or have considered it irrelevant [4]. The delay in getting good 3-D estimates comes from the fact that such estimates require several BO observations from different viewpoints. This makes landmark initialization in BO-SLAM difficult, to the point that satisfactory methods able to exploit all the geometrical information provided by the cameras have only recently become available. We have witnessed an evolution of the algorithms as follows. First, *delayed landmark initialization* methods attempted to obtain a full 3-D estimate before initialization via several observations from different viewpoints. Davison [3] showed real-time feasibility of monocular SLAM with affordable hardware, using the original extended Kalman filter (EKF) SLAM algorithm for all but the unmeasured landmark's depth, and a separate particle filter to estimate this depth. Initialization was *deferred* to the moment when the depth estimate was good enough. The consequence of a delayed scheme is that we can only initialize landmarks with enough parallax, i.e., those that are close to the camera and situated perpendicularly to its trajectory, and therefore, the need to operate in room-size scenarios with lateral motions. Second, Solà *et al.* [1] showed that *undelayed landmark initialization* (mapping the landmarks from their first, partial observation) was needed when considering low parallax landmarks, i.e., those that are remote and/or situated close to the motion axis. This permits mapping larger scenes while performing frontal trajectories. Third, Civera *et al.* [5] have recently achieved the mapping of landmarks up to infinity, due to an undelayed initialization via an *inverse depth parameterization* (IDP). IDP has also been developed by Eade *et al.* [6] in a FastSLAM2.0 context. Today, the monocular SLAM systems exploit the geometrical information in its entirety: from the first observation, independently of the sensor's trajectory, and up to the infinity range.

90 B. Structure From Motion (SFM)

91 Monocular SLAM compares to a similar problem solved
92 by the vision community: the structure from motion problem
93 (SFM). In SFM, the goal is to determine, from a collection of
94 images and up to an unrecoverable scale factor, the 3-D structure
95 of the perceived scene and all 6-D camera poses from where the
96 images were captured. When compared to SLAM, the structure
97 plays the role of the map, while the set of camera poses defines
98 all the successive observer's localizations.

99 Roboticians often claim that the main difference between
100 SFM and SLAM is that the former is solved offline via
101 the iterative nonlinear optimization method known as bundle
102 adjustment (BA) [7], while the latter must be incremen-
103 tally solved online, thus making use of stochastic estimators
104 or *filters* that naturally provide incremental operation. This
105 has been true for some years (today, SLAM is also solved
106 online with iterative optimization [8]), but does not tell the
107 whole story. The differences between SFM and SLAM are
108 not only in the methods but also in the objectives, meaning
109 that similar aspects of similar problems are given different
110 priorities.

111 In particular, SFM exploits the visual information in its en-
112 tirety without the difficulties encountered in monocular SLAM.
113 Let us try to understand this curious fact. SFM puts the struc-
114 ture as a final objective, i.e., as a result of the whole process,
115 and the emphasis is placed on minimizing the errors in the
116 *measurement space*, thus using all the measured information.
117 On the other hand, the SLAM map has a central role, with
118 some of the operations (and particularly landmark initializa-
119 tion) being performed in map space, which is the system's *state*
120 *space*. The fact that this state space is not statically observable,
121 because it is of higher dimension than the observation space,
122 leads to the difficulties exposed before. As an informal attempt
123 to fill this gap, we could say that modern undelayed methods
124 for monocular SLAM, with partial landmark initialization and
125 partial updates, are almost equivalent to an operation in the
126 measurement space: the information is initialized in the map
127 space *partially*, i.e., exactly as it comes from the measurement
128 space. A similar point of view over this concept can be found
129 in [9].

130 C. Stereovision SLAM

131 Stereovision SLAM has also received considerable attention.
132 The ability of a stereo assembly to directly and immediately pro-
133 vide 3-D landmark estimates allows us to use the best available
134 SLAM algorithms and rapidly obtain good results with little
135 effort in the conceptual parts. Such SLAM systems consider
136 the stereo assembly as being a single monolithic sensor, capa-
137 ble of gathering 3-D geometrical information from the robot's
138 surroundings, e.g. [10]. This fact, which appears perfectly rea-
139 sonable, is the main paradigm that this paper questions. By
140 considering two linked cameras as a single 3-D sensor, SLAM
141 is unable to face the following two issues.

142 1) *Limited 3-D Estimability Range*: While cameras are ca-
143 pable of sensing visible objects that are potentially at infinity,
144 a stereo rig provides only reasonably good 3-D estimates up

to a limited range, typically from 3 m to a few tens of meters 145
depending on the baseline. Because classical, nonmonocular 146
SLAM algorithms expect full 3-D estimates for landmark ini- 147
tialization (i.e., they are reasoned in the map space), information 148
belonging to only this limited region can be used for SLAM. 149
This is really a pity; it is like if, having our two eyes, we were 150
obliged to neglect everything outside a certain range from us, 151
what we could call "*walking inside dense fog*." Without remote 152
landmarks, it is easy to lose spacial references, to become disori- 153
ented, and finally, find ourselves lost. Therefore, stereovision, 154
as it is classically conceived, is a bad starting point for visual 155
SLAM. 156

2) *Mechanical Fragility*: If we aim at extending the 3-D 157
estimability range beyond these few tens of meters, we need 158
to increase the stereo baseline while keeping or improving the 159
overall sensor precision. This is obviously a contradiction: larger 160
assemblies are less precise when using the same mechanical 161
solutions. In order to maintain accuracy with a larger assembly, 162
we must use more complex structures that will be either heavier 163
or more expensive, if not both. The result for moderately large 164
baselines (>1 m) is a sensor that is very easily decalibrated, 165
and therefore, almost useless. Large rigs, however, are very 166
interesting in outdoor applications because they allow farther 167
objects to be positioned, thus making them contribute to the 168
observability of the overall scale factor. This is especially true 169
in aerial and underwater settings where, without nearby objects 170
to observe, a small stereo rig provides no significant gain with 171
respect to a single camera. Self-calibration can compensate for 172
the inherent lack of stability of large camera rigs. It also allows 173
multicamera platforms to start operation without undergoing a 174
previous calibration phase, making on-field system deployment 175
and maintenance easier. 176

To our knowledge, the only SLAM work that goes beyond the 177
current stereoparadigm (apart from our conference paper [11]) 178
is the one by Paz *et al.* [12], which uses a small-baseline, fully 179
calibrated stereo rig. Matched features presenting significant 180
disparity are initialized as classical Euclidean landmarks, while 181
those presenting low disparities are treated with the inverse 182
depth algorithm. 183

D. Visual Odometry (VO)

184 One could say that, in terms of methodology, visual odom- 185
etry (VO) is to stereovision SLAM what SFM is to monocular 186
SLAM. VO is conceived to obtain the robot's ego motion from 187
a sequence of stereo images [13]. Visual features are matched 188
across two or more pairs of stereo images taken during the robot 189
motion. An iterative minimization algorithm, usually based on 190
BA, is run to recover the stereo rig motion, which is then trans- 191
formed into robot motion. For this, the algorithm needs to re- 192
cover the structure of the 3-D points that correspond to the 193
matched features. This structure is not exploited for other tasks 194
and can be usually discarded. Remarkably, when the structure 195
is coded in the measurement space (u, v, d) , a disparity $d \rightarrow 0$ 196
allows points at infinity to be properly handled [14]. This is also 197
accomplished by using homogeneous coordinates [7]. VO must 198
work in real time because robot localization is needed online. 199

200 Advanced VO solutions achieve very low drift levels after long
 201 distances by making use of: 1) hardware-based image process-
 202 ing with real-time construction and querying of large feature
 203 databases [15]; 2) dense image information matching via planar
 204 homographies and the use of the quadrifocal tensor [16]; or 3)
 205 bundle adjusting the set of N recent key frames together with
 206 additional fusion with an inertial measurement unit (IMU) [14].

207 E. Sensor Fusion in SLAM

208 The fact of SLAM being solved by filters allows us to envision
 209 SLAM systems as sensor fusion engines. Let us highlight some
 210 of the assets of filtering in sensor fusion.

- 211 1) *Multisensor operation*: Any number of differing sensors
 212 can be operated together in a consistent framework.
- 213 2) *Sensors self-calibration*: Unknown biases, gains, and
 214 other sensor's parameters can be estimated provided that
 215 they are observable [17].
- 216 3) *Desynchronized operation*: The data rates of all these sen-
 217 sors do not need to be synchronized.
- 218 4) *Decentralized operation*: Advanced filter formulations
 219 such as those using channel filters [18] achieve a decen-
 220 tralized operation that should permit live connection and
 221 disconnection of sensors without the need for filter repro-
 222 gramming or reparameterization.

223 This paper explores the first three points for the case of mul-
 224 tiple cameras.

225 SLAM systems naturally fuse information from both propri-
 226 oceptive (odometry, GPS, and IMU) and exteroceptive (range
 227 scanners, sonar, and vision) sensors into the map. But our in-
 228 terest here is in fusing several exteroceptive sensors. We can
 229 distinguish two cases.

- 230 1) *Sensors of different kind*: When using differing sensors
 231 (e.g., laser plus vision), the main problem is in finding a
 232 map representation well adapted to the different kinds of
 233 sensory data (i.e., the data association problem).
- 234 2) *Sensors of the same kind*: The perceived information is of
 235 the same nature. This makes appearance-based matching
 236 possible, and therefore, makes map building easier. Nev-
 237 ertheless, most of such SLAM systems do not take advan-
 238 tage of fusion. Instead, the extrinsic parameters linking
 239 the sensors are calibrated offline, and the set of sensors
 240 is treated as a single supersensor. This is the case for
 241 two 180° range scanners simulating a 360° one, and for
 242 the previously mentioned stereo rig simulating a 3-D sen-
 243 sor. A sensor-fusion approach in these cases should natu-
 244 rally bring the aforementioned advantages to the SLAM
 245 system.

246 F. Multicamera SLAM and the Aim of This Paper

247 The key idea of this paper is very simple: by employing
 248 the SLAM filter as a fusion engine, we will be able to use
 249 any number of cameras in any configuration. And, by treat-
 250 ing them as BO sensors with the modern undelayed initializa-
 251 tion methods, we will extract the entire geometrical information
 252 provided by the images. The filter—not the sensor—will be re-

sponsible for making the 3-D properties of the perceived world
 arise.

Applications may vary from the simplest stereo system,
 through robots with several differing cameras (e.g., a panoramic
 one for localization and a perspective one looking forward
 for reactive navigation), to multirobot cooperative SLAM
 where BO observations from different robots are used to
 determine the 3-D locations of very distant landmarks. Al-
 though there certainly exist issues concerning multicamera
 management, the main ideas we want to convey may be
 demonstrated with systems of just two cameras. In this pa-
 per, we will illustrate two cases: first, the case of a robot
 equipped with a stereo rig, with its cameras being treated
 as two individual monocular sensors and second, two cam-
 eras moving independently and mapping together an outdoors
 scene.

This paper draws on previous work published in the confer-
 ence paper [11] and the author's Ph.D. thesis [19]. These two
 works use the federated information sharing algorithm (FIS)
 in [1] to initialize the landmarks, which has been surpassed by
 the inverse depth methods (IDP) [5]. The present paper takes
 and extends all this research by developing a better founded jus-
 tification (providing a wider scope to the proposed concepts), by
 improving on the implementation with the incorporation of IDP
 in the algorithms, and by extending the experimental validation
 to a cooperative monocular SLAM setup.

This paper is organized as follows. Section II presents the
 main ideas that will be exploited later and revises some back-
 ground material for monocular SLAM. Section III explains how
 to set up multicamera SLAM, an application for stereo benches
 with self-calibration, and an application for two collaborative
 cameras. Section IV presents the perception and map manage-
 ment techniques used. Sections V and VI show the experimen-
 tal results, and finally, Section VII gives conclusions and future
 directions.

II. 3-D ESTIMABILITY IN VISUAL SLAM

In this section, we present the ideas that support our approach
 to visual SLAM. We make use of the concept of estimability,
 which will help understand the abilities of vision for observing
 3-D structure in the presence of uncertainty. We clarify the key
 properties of undelayed initialization in monocular SLAM, and
 remark its importance in multicamera SLAM. We also remind
 the key aspects of IDP-SLAM.

A. Geometrical Approach to 3-D Estimability

We are interested in finding the shape and dimensions of the
 3-D-estimable region defined by two monocular views.

For this, we start with a couple of ideas to help understand-
 ing the concept of estimability used. When a new feature is
 detected in an image, the backprojection of its noisy-measured
 position defines a conic-shaped *pdf* for the landmark position,
 called *ray*, which extends to infinity (see Fig. 1). Let us con-
 sider two features extracted and matched from a pair of images,
 corresponding to the same landmark: their backprojections are
 two conic rays A and B that extend to infinity. The angular

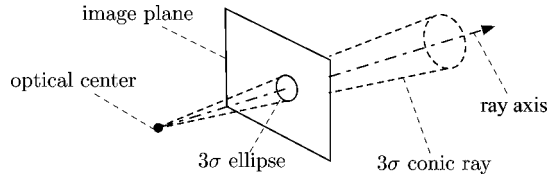


Fig. 1. Conic ray backprojects the elliptic representation of the Gaussian 2-D measure. It extends to infinity.

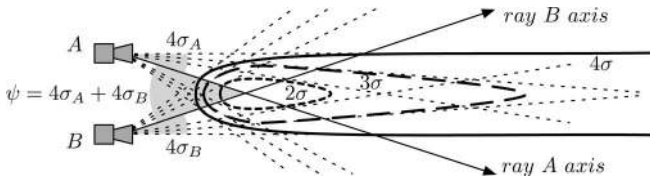


Fig. 2. Different regions of intersection for (solid) 4σ , (dashed) 3σ , and (dotted) 2σ ray widths when the outer 4σ bounds are parallel. (Shaded) The parallax or angle between rays axes A and B is $\psi = 4\sigma_A + 4\sigma_B$.

widths of these rays can be defined as a multiple of the standard deviations σ_A and σ_B of the angular errors (a composition of the cameras extrinsic and intrinsic parameters errors, and of the image processing algorithms accuracy). Informally speaking, we may say that the landmark's depth is fully estimated if the region of intersection of these rays is both *closed* and *sufficiently small*. If we consider, for example, the case where the two external 4σ bounds of the rays are parallel (see Fig. 2), then we can assure that the 3σ intersection region (which covers 98% probability) is *closed* and that the 2σ one (covering 74%) is *closed and small*. The ratio between the depth's standard deviation and its mean (a measure of linearity in monocular EKF-SLAM [1], [3]) is then better than 0.25. The *parallax* angle ψ between the two rays axes is therefore $\psi = 4(\sigma_A + \sigma_B) = \text{constant}$. This is the minimum parallax for full estimability.

In 2-D, we can plot the locus of constant estimability. In the case, where σ_A and σ_B can be considered constant, ψ is constant too, and from the inscribed angle theorem, the locus is then circular (Fig. 3, see also [19]). Landmarks inside this circle are considered *fully estimable*—and *partially* outside. In 3-D, the *fully 3-D estimable* region is obtained by revolution of this circle around the axis joining both cameras, producing a torus-shaped region with a degenerated central hole. This shape admits the following interpretations.

- 1) In a stereo configuration or for a lateral motion of a moving camera (see Fig. 3, left), the estimable region is located in front of the sensor. Beyond the region's border stereo provides no profit: if we want to consider distant landmarks, we have to use undelayed monocular techniques.
- 2) Depth recovery is impossible in the motion axis of a single camera moving forward (Fig. 3, right). Close to this axis, estimability is possible only if the region's radius becomes very large. This implies the necessity of very large displacements of the camera during the initializa-

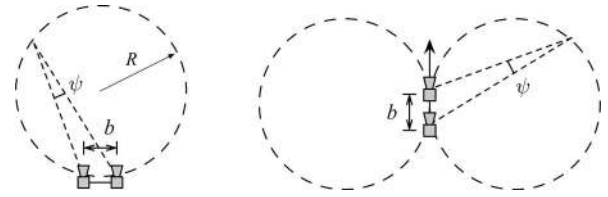


Fig. 3. Simplified depth estimability regions in a (left) stereo rig and (right) a camera traveling forward. The angle ψ is the one that assures estimability via triangulation from different viewpoints. The maximum range is $2R = b/\sin(\psi/2)$.

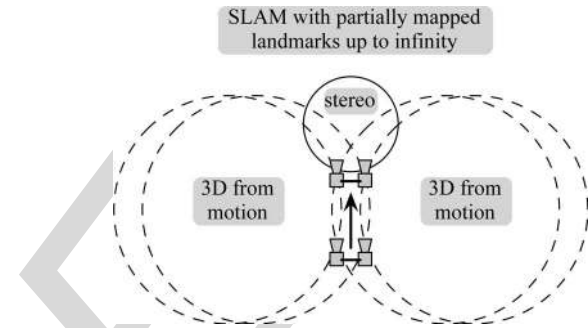


Fig. 4. Simplified depth estimability for a stereo rig moving forward. On both sides, estimability depends on the baseline gained by motion. In front, by stereo. Out of these bounds and up to infinity, landmarks are mapped partially. SLAM keeps incorporating the visual information due to the undelayed monocular methods, i.e., IDP in our case.

tion process. Again, this can be accomplished only with undelayed initializations.

- 3) By combining both monocular and stereovision, we get an instant estimability of close frontal objects while still utilizing the information of distant ones (see Fig. 4). Landmarks lying outside the estimability regions are not 3-D-estimable but, when initialized using undelayed monocular methods, they will contribute to constrain the camera orientation. Ideally, long-term observations of stable distant landmarks would completely cancel orientation drift (visual compass).

B. Monocular IDP-SLAM

The core algorithm of this paper is an EKF-SLAM with an IDP of landmarks during the initialization phase, as described in [5]. In IDP-SLAM, partially observed landmarks are coded as a 6-D-vector,

$$\mathbf{i} = [\mathbf{x}_0, \theta, \psi, \rho] \quad (1)$$

where \mathbf{x}_0 is the 3-D position of the camera at initialization time, (θ, ψ) are the elevation and azimuth angles in global frame defining the direction of the landmark's ray, and ρ is the inverse of the Euclidean distance from \mathbf{x}_0 to the landmark's position (notice that ρ is usually known as *inverse depth* but it is rather an inverse distance). After the first observation, all parameters of \mathbf{i} except ρ are immediately observable, and their values and covariances are obtained by proper inversion and linearization of the observation functions. The inverse depth ρ is initialized

369 with a Gaussian $\mathcal{N}(\rho - \bar{\rho}; \sigma_\rho^2)$ such that in the depth dimension
370 $s = 1/\rho$, we have

$$s_{(-n\sigma)} = \frac{1}{\bar{\rho} - n\sigma_\rho} = \infty \quad (2)$$

$$s_{(+n\sigma)} = \frac{1}{\bar{\rho} + n\sigma_\rho} = s_{\min} \quad (3)$$

371 with s_{\min} the minimum considered depth and n the inverse depth
372 shape factor. This gives $\bar{\rho} = 1/(2s_{\min})$ and, more remarkably

$$n\sigma_\rho = \bar{\rho}. \quad (4)$$

373 Importantly, values of $1 \leq n \leq 2$ assure from (2) that the infinity
374 range is included in the parametrization with ample probability.

375 On subsequent updates, IDP achieves correct EKF operation
376 (i.e., quasi-linear behavior) along the whole ray as long as the
377 parallax shown by the new viewpoint is not too large. The linear-
378 earity test in [20] is regularly evaluated. If passed, the landmark
379 can be safely transformed into a 3-D Euclidean parametrization.

380 III. MULTICAMERA SLAM

381 The general scheme for the multicamera SLAM system is
382 presented in this section. This scheme is particularized to deal
383 with two different problems. The first one is the automatic self-
384 calibration of a stereo rig while performing SLAM. The second
385 one is a master-lave solution to cooperative monocular SLAM.
386 Both setups are explained here, and their corresponding experi-
387 ments are presented in Sections V and VI.

388 A. System Overview

389 We implement the multicamera SLAM system as follows. A
390 central EKF-SLAM will hold the stochastic representation of
391 the set of all cameras \mathcal{C}_i plus the set of landmarks \mathcal{L}_j

$$X^\top = [\mathcal{C}_1^\top \ \cdots \ \mathcal{C}_N^\top \ \mathcal{L}_1^\top \ \cdots \ \mathcal{L}_M^\top] \quad (5)$$

392 where the cameras states contain position and orientation quaternion
393 $[\mathcal{C}_i = (\mathbf{r}_i, \mathbf{q}_i) \in \mathbb{R}^7]$, and landmarks can be coded either
394 in inverse depth ($\mathcal{L}_j = \mathbf{i}_j \in \mathbb{R}^6$) or in Euclidean coordinates
395 ($\mathcal{L}_j = \mathbf{p}_j \in \mathbb{R}^3$). Any number of cameras can be considered
396 this way. As each camera needs to remain localized properly,
397 it needs to observe a minimum number of landmarks at each
398 frame. The algorithm's complexity increases linearly with the
399 number of cameras if this number is small with respect to the
400 map.

401 For camera motions, we consider two possible models. In
402 the first one, a simple odometer provides motion predictions
403 $[\Delta x, \Delta y, \Delta \psi]$ in the robot's local 2-D plane. Gaussian uncer-
404 tainties are added to the 6-DOF linear and angular components
405 $[x, y, z, \phi, \theta, \psi]$ with a variance proportional to the measured
406 forward motion Δx

$$\{\sigma_x^2, \sigma_y^2, \sigma_z^2\} = k_L^2 \cdot \Delta x \quad (6)$$

$$\{\sigma_\phi^2, \sigma_\theta^2, \sigma_\psi^2\} = k_A^2 \cdot \Delta x. \quad (7)$$

407 The variance in $[\phi, \theta, \psi]$ is mapped to the quaternion space using
408 the corresponding Jacobians.

The second model is a 6-DOF constant velocity model

$$\mathbf{r}^+ = \mathbf{r} + \mathbf{v} \Delta t$$

$$\mathbf{q}^+ = \mathbf{q} \times v2q(\boldsymbol{\omega} \Delta t)$$

$$\mathbf{v}^+ = \mathbf{v} + \boldsymbol{\eta}_v$$

$$\boldsymbol{\omega}^+ = \boldsymbol{\omega} + \boldsymbol{\eta}_\omega$$

410 where $()^+$ means the updated value, \times is the quaternions prod-
411 uct, and $v2q(\boldsymbol{\omega} \Delta t)$ transforms the local incremental rotation
412 vector $\boldsymbol{\omega} \Delta t$ into a quaternion (quaternions are systematically
413 normalized). This way, the camera state vector \mathcal{C}_i is augmented
414 to $\mathcal{C}_i = (\mathbf{r}_i, \mathbf{q}_i, \mathbf{v}_i, \boldsymbol{\omega}_i) \in \mathbb{R}^{13}$. At each time step, perturbations
415 $\{\boldsymbol{\eta}_v, \boldsymbol{\eta}_\omega\} \sim \mathcal{N}(0; \{\sigma_v^2, \sigma_\omega^2\})$ add variances to the linear and
416 angular velocities proportionally to the elapsed time Δt

$$\sigma_v^2 = k_v^2 \cdot \Delta t \quad (8)$$

$$\sigma_\omega^2 = k_\omega^2 \cdot \Delta t. \quad (9)$$

417 The events of camera motion, landmark initialization, and
418 landmark observation are handled as in regular IDP-SLAM by
419 just selecting the appropriate block elements from the SLAM
420 state vector and covariances matrix, and applying the corre-
421 sponding motion or observation models. For example, at the
422 observation of landmark j from camera i , we would use the
423 function $\mathbf{u}_j^i = \mathbf{h}(\mathcal{C}_i, \mathcal{L}_j)$, which will be explained later for the
424 case of an IDP ray [see 11]. Before transforming IDP rays into
425 points, the linearity test in [20] needs to hold for all cameras.

426 B. Stereo SLAM With Extrinsic Self-Calibration

427 Our approach is relevant to fully calibrated stereo rigs if they
428 are small (10–20 cm, as in [12]) or if, having long baselines, their
429 main extrinsic parameters can be continuously self-calibrated.

430 Not all of the six extrinsic parameters of a stereo rig (three for
431 translation, three for orientation) need to be calibrated. In fact,
432 the notion of *self-calibration* inherently requires the system to
433 possess its own gauge. In our case, the metric dimensions or
434 *scale factor* of the whole world-robot system can only be ob-
435 tained either from the stereo rig baseline, which is one of the
436 extrinsic parameters (then, it makes no sense to self-calibrate
437 the gauge), or from odometry, which is often much less accurate
438 than any coarse measurement we could make of this baseline.
439 Additionally, as cameras are actually angular sensors, vision
440 measurements are much more sensitive to the cameras orienta-
441 tions than to any translation parameter. This means that vision
442 measurements will contain little information about these trans-
443 lation parameters. In consequence, self-calibration may concern
444 only orientation, and more precisely, the orientation of one cam-
445 era with respect to the other. The error of the reconstructed map's
446 scale factor will be the same as the relative error of the baseline
447 measurement.

448 With these assumptions, our self-calibration solution is
449 straightforward: for the second camera, we just include its ori-
450 entation in the map and let EKF make the rest. The state vector
451 (5) is modified and written as

$$X^\top = [\mathcal{R}^\top \ \mathbf{q}_R^\top \ \mathcal{L}_1^\top \ \cdots \ \mathcal{L}_M^\top]$$

452 where \mathcal{R} and $\mathcal{L}_1 \cdots \mathcal{L}_M$ are the robot pose and landmarks map.
 453 The left camera pose \mathcal{C}_L has a fixed transformation with respect
 454 to the robot, and \mathbf{q}_R is the orientation part of the right-hand
 455 camera \mathcal{C}_R in the robot frame. The time-evolution function of
 456 the angular extrinsic parameters is simply $\mathbf{q}_R^+ = \mathbf{q}_R + \gamma$, where
 457 γ is a white, Gaussian, low-energy process noise that accounts
 458 for eventual decalibrations, e.g., due to vibrations. For short-
 459 duration experiments, we set $\gamma = 0$. A coarse analysis of the
 460 stereo structure's mechanical precision will be enough to set the
 461 initial uncertainty to a value of the order of 1° or 2° per axis.
 462 This can be reduced to a few tenths of degree in cases where we
 463 dispose of previous calibrated values about which we are not
 464 confident anymore.

465 C. Cooperative Multicamera SLAM

466 The ideal, most general case of cooperative SLAM (5), corre-
 467 sponds to a (not too large) number of cameras moving indepen-
 468 dently. Each camera is able to manage its own measurements
 469 and communicates directly with the map. The aim of this com-
 470 munication is to obtain information about existing landmarks
 471 to get localized, and provide information about new or reob-
 472 served landmarks. This way, the algorithms to be executed by
 473 each camera are absolutely symmetrical, without any kind of
 474 hierarchy. A simplified implementation considers cameras with
 475 different privileges.

476 In our particular case, the cooperative SLAM system consid-
 477 ers two cameras. One of them takes the role of *master*, and
 478 is responsible for all landmarks detection and initialization.
 479 The second one acts as the *slave*. It follows the master at a
 480 close distance and reobserves the SLAM map that is being
 481 built by the master. By doing so, it provides a second view-
 482 point to landmarks just initialized, accelerating the convergence
 483 of the map. The master and slave trajectories are highly in-
 484 dependent, and for instance, they can cross paths. The only
 485 requirement is to look in the same direction. A trivial exten-
 486 sion to more than two cameras consists in including additional
 487 slaves.

488 IV. PERCEPTION AND MAP MANAGEMENT

489 Active search (AS, nicely described in [21] and also referred
 490 to as *top-down* in [6]) is a powerful framework for real-time
 491 image processing within SLAM. It has been successfully used in
 492 several monocular SLAM works [3], [5], [11], using a diversity
 493 of techniques for landmark initialization. The idea of AS is to
 494 exploit the information contained in the map to predict a number
 495 of characteristics of the landmarks to observe. AS is helpful in
 496 solving the following issues:

- 497 1) selecting interesting image regions for initialization;
- 498 2) selecting the most informative landmarks to measure;
- 499 3) predicting where in the image they may be found, and with
 500 which probability;
- 501 4) predicting the current landmark's appearance to maximize
 502 the chances of a successful match.

A. Feature Detection and Initialization

503

504 Based on the projection of the map information into the master
 505 image, a heuristic strategy is used to select a region of interest
 506 for a new initialization: we divide the image with a grid and
 507 randomly select a grid element with no landmarks inside. We
 508 extract the strongest Harris point [22] in this region and validate
 509 it if its strength is above a predefined threshold. We store a small
 510 rectangular region or *patch* of 15×15 pixels around the point
 511 as the landmark's appearance descriptor, together with the pose
 512 of the camera. Finally, we initialize the IDP ray in the SLAM
 513 map.

B. Expectations: The Active Search Regions

514

515 Some considerations about AS can be made for its usage in
 516 multicamera IDP-SLAM to improve performance. We use for
 517 this the \mathcal{E}_1 and \mathcal{E}_∞ ellipses, defined and explained as follows.

518 1) \mathcal{E}_1 *Ellipse: Expectation of the Inverse Depth Ray*: The
 519 inverse depth ray (1) is easily projected into a camera. We take
 520 the transformation to camera frame given in [5]:

$$\mathbf{h}_1^c = \mathbf{R}(\mathbf{q})^\top (\rho(\mathbf{x}_0 - \mathbf{r}) + \mathbf{m}(\theta, \psi)) \quad (10)$$

521 where $\mathbf{R}(\cdot)$ is the rotation matrix corresponding to the camera
 522 orientation \mathbf{q} and \mathbf{r} is the current camera position. This value
 523 is then projected into the camera, described by intrinsic and
 524 distortion parameters \mathbf{k} and \mathbf{d} (we use a classical radial distortion
 525 model of up to three parameters, which is inverted as explained
 526 in [19]). Let us call $\mathcal{K} = (\mathbf{k}, \mathbf{d})$ the camera parameters, $\mathcal{C} =$
 527 (\mathbf{r}, \mathbf{q}) the camera pose, and $\mathbf{i} = (\mathbf{x}_0, \theta, \psi, \rho)$ the IDP ray. The
 528 observation function is

$$\mathbf{u} = \mathbf{h}_1(\mathcal{C}, \mathcal{K}, \mathbf{i}) + \eta = \text{project}(\mathbf{h}_1^c, \mathcal{K}) + \eta \quad (11)$$

529 where $\text{project}(\cdot)$ takes into account the camera model (we use
 530 perspective cameras) and η is the pixel Gaussian noise, with
 531 covariance \mathbf{R} .

532 We define the \mathcal{E}_1 ellipse as the Gaussian expectation
 533 $\mathcal{E}_1(\mathbf{u}) \triangleq \mathcal{N}(\mathbf{u} - \bar{\mathbf{e}}_1; \mathbf{E}_1)$, with \mathbf{u} being the pixel position, and
 534 with mean and covariances matrix

$$\bar{\mathbf{e}}_1 = \mathbf{h}_1(\bar{\mathcal{C}}, \mathcal{K}, \bar{\mathbf{i}}) \quad (12)$$

$$\mathbf{E}_1 = [\mathbf{H}_c \mathbf{H}_i] \mathbf{P}_{c,i} [\mathbf{H}_c \mathbf{H}_i]^\top + \mathbf{R}. \quad (13)$$

535 Here, \mathbf{H}_c and \mathbf{H}_i are the Jacobians of \mathbf{h}_1 with respect to the
 536 uncertain parameters \mathcal{C} and \mathbf{i} , \bullet are variable estimates from
 537 the SLAM map, and $\mathbf{P}_{c,i}$ is the joint covariances matrix (all
 538 correlations and cross correlations) of \mathcal{C} and \mathbf{i} , also from the
 539 map. In AS, \mathcal{E}_1 is usually gated at 3σ , giving place to an elliptic
 540 region in the image where the landmark must project with 98%
 541 probability. However, this is not necessarily true in cases of
 542 noticeable parallax, as we examine now.

543 At landmark initialization, its inverse depth ρ is initialized
 544 according to (2)–(4). When considering 3σ uncertainty regions,
 545 (4) implies that ρ can go negative with a nonnegligible probab-
 546 ility, meaning that the coded landmarks might be situated *behind*
 547 *the camera*. This becomes evident when projecting the IDP ray
 548 into a second camera presenting some parallax: the projected
 549 3σ \mathcal{E}_1 ellipse contains a region with negative disparity (see

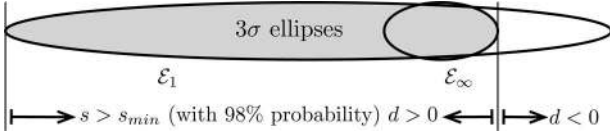


Fig. 5. 3σ search region defined by the \mathcal{E}_1 ellipse contains a significant part that corresponds to negative disparities $d < 0$, where the feature should not be searched. The final 3σ search region (gray) is defined by the \mathcal{E}_1 and \mathcal{E}_∞ ellipses. The rightmost 3σ border of \mathcal{E}_∞ is where the probability to find the projection of the infinity point has fallen below 2%.

Q1

550 Fig. 5). It is desirable to limit the search area to values of only
 551 positive disparity for two reasons: the correlation-based search
 552 (one of the most time-consuming processes) is faster and the
 553 possibility of including false matches as outliers is diminished.
 554 With nonrectified images and/or camera sets with uncertain ex-
 555 trinsic parameters, determining the null disparity bound is not
 556 straightforward. One solution is to use the \mathcal{E}_∞ ellipse, which we
 557 introduce in the following paragraph.

558 2) \mathcal{E}_∞ Ellipse: *Expectation of the Infinity Point*: The infinity
 559 point is easily projected by considering the transformation (10)
 560 with $\rho \rightarrow 0$

$$\mathbf{h}_\infty^c \approx \mathbf{R}(\mathbf{q})^\top \mathbf{m}(\theta, \psi) \quad (14)$$

561 where only the camera orientation \mathbf{q} and the ray's direction
 562 angles (θ, ψ) are present (the visual compass). Proceeding as
 563 before, we obtain the definition of the ellipse $\mathcal{E}_\infty(\mathbf{u}) \triangleq \mathcal{N}(\mathbf{u} -$
 564 $\bar{\mathbf{e}}_\infty; \mathbf{E}_\infty)$ as

$$\bar{\mathbf{e}}_\infty = \mathbf{h}(\bar{\mathbf{q}}, \bar{\mathcal{K}}, \bar{\theta}, \bar{\psi}) \quad (15)$$

$$\mathbf{E}_\infty = [\mathbf{H}_q \mathbf{H}_\theta \mathbf{H}_\psi] \mathbf{P}_{\{q, \theta, \psi\}} [\mathbf{H}_q \mathbf{H}_\theta \mathbf{H}_\psi]^\top + \mathbf{R} \quad (16)$$

565 where $\mathbf{P}_{\{q, \theta, \psi\}}$ is the joint covariances matrix of the uncertain
 566 parameters. The \mathcal{E}_∞ 3σ region is composed of the previous \mathcal{E}_1
 567 region, as indicated in Fig. 5, to define the search area.

568 C. Selection of the Best Map Updates

569 Following the AS approach in [23], a predefined number of
 570 landmarks with the biggest \mathcal{E}_1 ellipse surfaces are selected in
 571 each camera as those being the most interesting to be measured.
 572 For each camera, we organize all candidates (visible landmarks)
 573 in descending order of expectation surfaces, without caring if
 574 they are points or rays. We update at each frame a predefined
 575 number of them (usually around 10, and no more than 20).
 576 Updates are processed sequentially, with all Jacobians being
 577 recalculated each time to minimize the effect of linearization
 578 errors.

579 D. Feature Matching: Affine Patch Warping

580 AS continues by *warping* the stored patch and searching for
 581 a correlation peak inside the search area earlier. The objec-
 582 tive of warping is to predict the landmark's current appearance,
 583 maximizing the chances for a good match. In the absence of dis-
 584 tortion, a planar homography $\mathbf{H} \in \mathbb{R}^{3 \times 3}$, defined in the homo-
 585 geneous spaces, would be desirable [24]. This type of warping
 586 requires the online estimation of the patch normal in the 3-D

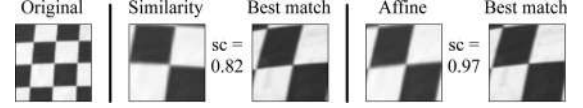


Fig. 6. Similarity and affine warping on a sample patch. From left to right: original patch; similarity warped patch ($\sim 180\%$ scale, 10° rotation); best match in a later image affected by distortion and its zero mean normalized cross correlation (ZNCC) score (0.82); affine warped patch; best match and score (0.97). The affine warping contains a significant skew component mainly due to image distortion. The improvement in the ZNCC score is very important.

587 space, and may become very time-consuming. A good simplifi-
 588 cation considers this normal fixed at the initial visual axis [23].
 589 Further simplification applies just a similarity transformation
 590 $\mathbf{T} = s\mathbf{R} \in \mathbb{R}^{2 \times 2}$ in the image Euclidean plane [19]. This ac-
 591 counts only for scale changes s and rotations \mathbf{R} obtained from
 592 the stored information (landmark position, camera initial, and
 593 current poses). However, in the presence of distortion, features
 594 lying close to the image borders suffer from additional deforma-
 595 tions. We developed a warping approach that easily adds a
 596 skew component to the operator \mathbf{T} (thus achieving fully affine
 597 warping, but not perspective warping; Fig. 6), based on the Ja-
 598 cobian of the function linking the first observation to the current
 599 one. Let us consider the backward observation model $\mathbf{g}(\cdot)$ for a
 600 camera A at initialization time $t = 0$, and the observation model
 601 $\mathbf{h}(\cdot)$ for a different camera B at current time $t \geq 0$

$$\mathbf{p} = \mathbf{g}(\mathcal{C}_A(0), \mathcal{K}_A, \mathbf{u}_A(0), s_A)$$

$$\mathbf{u}_B(t) = \mathbf{h}(\mathcal{C}_B(t), \mathcal{K}_B, \mathbf{p}).$$

602 Here, \mathbf{p} is the landmark's position, $\mathcal{K}_i = (\mathbf{k}_i, \mathbf{d}_i)$ are the intrin-
 603 sic and distortion parameters of camera i , $\mathbf{u}_i(t)$ is the measured
 604 pixel, and s_A is the landmark's depth with respect to the initial
 605 camera. We can compose these functions to obtain the expres-
 606 sion linking the initial and the current pixels
 607

$$\mathbf{u}_B(t) = \mathbf{h}[\mathcal{C}_B(t), \mathcal{K}_B, \mathbf{g}(\mathcal{C}_A(0), \mathcal{K}_A, \mathbf{u}_A(0), s_A)]. \quad (17)$$

608 When all but the pixel positions are fixed, this represents an
 609 invertible mapping $\mathbb{R}^2 \mapsto \mathbb{R}^2$ from the pixels in the first image
 610 to the pixels in the current one. The local linearization around
 611 the initially measured pixel defines an affine warping expressed
 612 by the Jacobian matrix

$$\mathbf{T} = \left. \frac{\partial \mathbf{u}_B}{\partial \mathbf{u}_A} \right|_{(\mathcal{C}_A(0), \mathcal{C}_B(t), \mathcal{K}_A, \mathcal{K}_B, \mathbf{u}_A(0), s_A)}. \quad (18)$$

613 By defining $\tilde{\mathbf{u}}_i$ as the coordinates of the patch in camera i , with
 614 the central pixel \mathbf{u}_i as the origin, we have $\tilde{\mathbf{u}}_B(t) = \mathbf{T} \tilde{\mathbf{u}}_A(0)$.
 615 Based on this mapping, we use linear interpolation of the pixels'
 616 luminosity to construct the warped patch.

617 V. EXPERIMENT 1: STEREO SLAM WITH SELF-CALIBRATION

618 The "White-board" indoor experiment aims at demonstrating
 619 stereovision SLAM with self-calibration. A robot with a stereo
 620 head looking forward is run for about 10 m in straight line inside
 621 the robotics laboratory at the LAAS (see Fig. 7). Over 500 image
 622 pairs are taken at approximately 5-Hz frequency. The robot

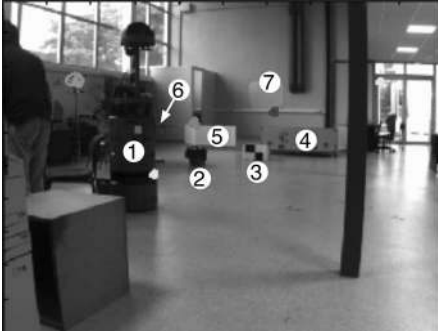


Fig. 7. Laboratoire d'Analyse et d'Architecture des System (LAAS) robotics laboratory. The robot will approach the scene in a straightforward trajectory. We notice in the scene the presence of a robot ①, a bin ②, a box ③, a trunk ④, a fence ⑤, a table ⑥ (hidden by the robot in this image), and the white board ⑦ at the end wall.

TABLE I
STEREO RIG PARAMETERS IN THE "WHITE-BOARD" EXPERIMENT

Scope	Parameters = Values
Dimensions	Baseline = 33 cm
Orientation - Euler	$\{\phi, \theta, \psi\} = \{0^\circ, 5^\circ, 0^\circ\}$
Cameras	$\{\text{resolution}, \text{FOV}\} = \{512 \times 384 \text{ pix}, 55^\circ\}$
Right camera uncertainties	$\{\sigma_\phi, \sigma_\theta, \sigma_\psi\} = \{1^\circ, 1^\circ, 1^\circ\}$

623 moves towards the objects to be mapped at 0.15 m/s. The stereo
624 rig consists of two intrinsically calibrated cameras arranged
625 as indicated in Table I. The orientations of both cameras are
626 specified with respect to the robot frame. The left camera is taken
627 as reference, thus deterministically specified, and the orientation
628 of the right one is initialized with an uncertainty of 1° standard
629 deviation. We use the odometry model (Section III-A) with
630 $k_L = 0.1 \text{ m}/\sqrt{\text{m}}$ and $k_A = 0.05 \text{ rad}/\sqrt{\text{m}}$.

631 We show details and results on the self-calibration procedure
632 and the metric accuracy of the resulting map. The mapping
633 process can be appreciated in the movie `whiteboard.mov` in
634 the multimedia section.

635 A. Self-Calibration

636 We plot in Fig. 8 left the evolution of the three self-calibrated
637 angles. We have also used the shape of the \mathcal{E}_∞ ellipses to pro-
638 vide additional qualitative evidence of the calibration process
639 (Fig. 9 and movie `whiteboard - e1nf.mov`). We observe the
640 following behavior.

641 1) *Pitch* θ : The pitch angle (cameras tilt, 5° nominal value) is
642 observable from the first matched landmark. It rapidly converges
643 to an angle of 4.77° and remains very stable during the whole
644 experiment.

645 2) *Roll* ϕ : Roll angle is observable after at least two land-
646 marks are observed from the right camera. Once this condition
647 holds, convergence occurs relatively fast.

648 3) *Yaw* ψ : Yaw angle is very weakly observable because
649 it is coupled with the landmarks depths: both yaw angle and
650 landmark depth variations produce a similar uncertainty growth
651 in the right image. For this reason, yaw converges slowly, only
652 showing reasonable convergence after some 50 frames.

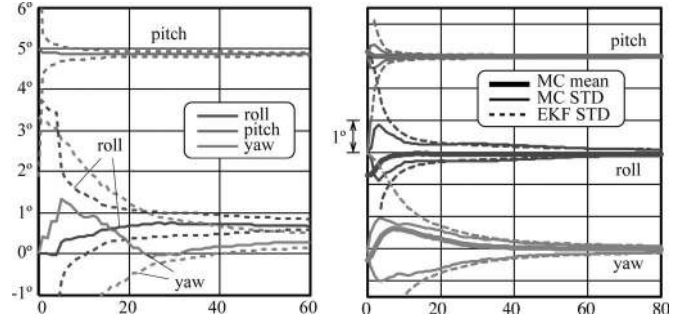


Fig. 8. Extrinsic self-calibration. (Left) The three Euler angles of the right camera orientation with respect to the robot during the first 60 frames. The 3σ bounds are plotted in dotted line showing consistent estimation. (Right) Error analysis after 100 MC runs using 200 frames per run (only the first 80 frames are shown). The thick solid lines represent the means over the 100 runs. The 3σ bounds for each angle are plotted using thin solid lines. The dotted lines represent the averaged 3σ bounds estimated by the EKF, showing consistent calibration.

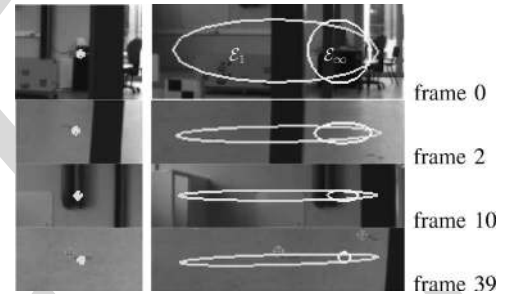


Fig. 9. Evolution of the \mathcal{E}_1 and \mathcal{E}_∞ ellipses during calibration. On the left column, newly detected pixels in the left image. On the right, expectations in the right image (green) \mathcal{E}_1 and (yellow) \mathcal{E}_∞ of the newly initialized IDP rays (i.e., still with the full initial uncertainty in ρ). At frame 0, initial uncertainties of 1° result in a big, round \mathcal{E}_∞ ellipse. After the first updated landmark from the left camera (frame 2), the uncertainty in pitch decreases and \mathcal{E}_∞ becomes flat. Successive updates further refine the calibrated angles. The yaw angle takes longer to converge, but the tiny \mathcal{E}_∞ in frame 39 shows that the calibration is already finished. The portion of the green ellipse on the right side of the yellow one corresponds to negative disparities and is not searched for matches. This portion is larger as parallax increases.

TABLE II
MC ANALYSIS OF THE SELF-CALIBRATION

Angle	MC mean	STD	EKF STD	Offline	STD
roll	0.69°	0.028°	0.018°	0.61°	0.013°
pitch	4.77°	0.003°	0.005°	4.74°	0.099°
yaw	0.33°	0.021°	0.016°	0.51°	0.109°

In Fig. 8 right, we plot results of a Monte Carlo (MC) analysis, run over the data of this experiment, for the mean and standard deviation of the Euler angles of the right camera. Because all MC runs are extracted from the same sequence, we tried to maximize their independence by using a different *random seed* in the algorithm (acting in the random selection of the initialization region, Section IV-A), and by starting each run at a different frame. The figure shows that the dynamic estimation is consistent (the EKF estimated sigmas are larger than the MC ones). After 200 frames, we compare these values with those of the offline calibration [25]. Table II summarizes these results, showing MC [(means and standard deviations (STD))] and Kalman Filter (EKF, showing the estimated STD). All

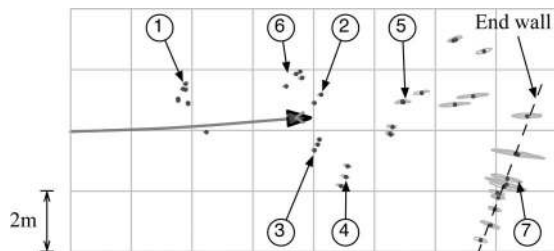


Fig. 10. Map produced during the “white board” experiment. We marked the mapped robot ①, the bin ②, the box ③, the trunk ④, the fence ⑤, the table ⑥, and the white board ⑦ at the end wall.

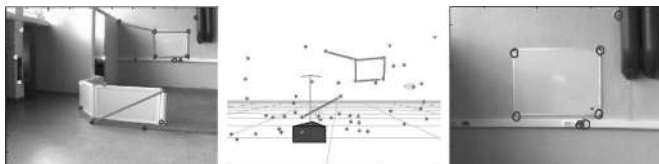


Fig. 11. Metric mapping. The magnitudes of some segments in the real laboratory are compared to those in the map (red lines). Ground truth corresponds to metric measurements of the distances between landmarks that are identified by zooming in the last image of the experiment (right) and translated to the real world. Thirteen points on the end wall are tested for coplanarity.

TABLE III
WHITE BOARD: MAP TO GROUND TRUTH TOMPARISON

segment	board	board	board	board	wall	fence
real (cm)	116	86	117	88	136	124
mapped	116.6	87.2	115.8	87.0	135.1	125.5
STD	0.91	0.81	1.21	0.52	1.06	1.32

666 self-calibrated values lie within the 3σ bounds defined by the
667 offline mean and STD values.

668 B. Metric Accuracy

669 We show in Fig. 10 a top view of the map generated during
670 this experiment. To contrast this map against reality, two tests
671 are performed: planarity and metric scale (see Fig. 11): 1) the
672 four corners of the white board are taken together with nine
673 other points at the end wall to test coplanarity: the 13 mapped
674 points are found to be coplanar within 4.9 cm STD; 2) the
675 lengths of the real and mapped segments marked in Fig. 11
676 are summarized in Table III. The white board has a physical
677 size of 120 cm \times 90 cm, but we take real measurements from
678 the approximated corners where the features are detected. We
679 observe errors in the order of 1 cm for landmarks that are still
680 about 4 m away from the robot.

681 VI. EXPERIMENT 2: COOPERATIVE MONOCULAR SLAM

682 This experiment shows independent cameras collaborating to
683 build a 3-D map using exclusively bearings-only observations.
684 Two independent cameras are placed on top of two bicycles
685 looking forward, moving on different trajectories in the park-
686 ing of the LAAS (see Fig. 12). Over 1000 images are taken
687 by each camera at 15-Hz frequency, 512 \times 384 pixel resolution,
688 100° field of view (FOV), and are processed offline. The cam-



Fig. 12. Snapshots of master and slave sequences in cooperative SLAM. Faraway landmarks (e.g., black arrowed), still initialized as rays (red), are the ones fixing the orientation. Nearby landmarks, usually as Euclidean points (blue), maintain the metric. A virtual model of the master camera is visible from the slave camera (white arrowed). See `cooperativeSLAM.mov`.

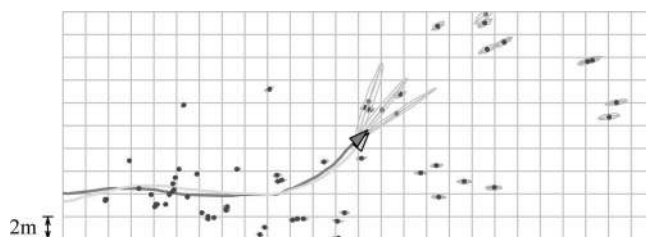


Fig. 13. Top view of the map produced by cooperative SLAM of two independent cameras, and their crossing trajectories. The grid spacing is 2 m.

eras travel approximately 28 m observing landmarks beyond 689
690 60 m. As in the previous experiment, the left camera is the master.
691 The two trajectories start parallel to each other, separated
692 75 cm perpendicularly to the motion direction. The reference
693 frame is defined by the master camera initial position and
694 orientation, which are initialized with null uncertainty. The scale
695 factor is determined by the initial baseline of 75 cm, meaning
696 that the position of the slave camera in the lateral Y -axis is also
697 initialized with null uncertainty. The orientations of the slave
698 camera start with an uncertainty of 2° STD, and its position in
699 the frontal Y - and vertical Z -axes with $75 \text{ cm} \cdot \sin(2^\circ) = 2.6 \text{ cm}$
700 STD. With these uncertainties, the experiment’s initial configura-
701 tion can be set up manually by just observing the images and
702 centering the projections of some distant object. We use two
703 independent constant-velocity models with $k_v = 0.3 \text{ m/s} \cdot \sqrt{s}$
704 and $k_w = 0.3 \text{ rad/s} \cdot \sqrt{s}$. The measurement noise is 1 pixel.

Landmarks at infinity, illumination changes and few salient 705
706 features are some characteristics of this outdoors scene. It
707 presents relatively few stable landmarks, something that makes
708 the correct operation of the SLAM system difficult. In the case
709 of having crossing trajectories, the problem of one camera oc-
710 cluding the other could appear and severely affect the image
711 processing. To avoid this, we decided to take both image se-
712 quences shifted in time, i.e., one after the other, and make them
713 overlap for processing. The mapping process is presented in
714 the movie `cooperativeSLAM.mov` in the multimedia section.
715 Fig. 13 shows the top view of the map and the camera trajec-
716 tories generated during this experiment.

A proper metrical evaluation of this experiment is difficult; 717
718 having a variable baseline does not allow us to compare the re-
719 sults, because there is no knowledge of the ground truth. In order
720 to evaluate this approach, we consider the setup in experiment

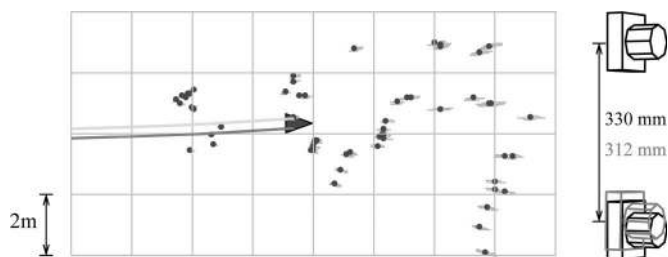


Fig. 14. Final map in the “white board” setup using the cooperative monocular SLAM algorithm. The cameras are modeled as being entirely independent using the same data and initial configuration as in Experiment 1. The stereo rig on the right shows (red) the final estimated relative position compared with (black) ground truth.

721 1 and apply the same algorithm. The new experiment consists
 722 of recovering the full extrinsic calibration, which is fixed in re-
 723 ality, considering both cameras as independent. Again, we use
 724 a constant-velocity model for each camera. The initial setup
 725 including uncertainties is as in experiment 1.

726 Fig. 14 shows the obtained map. We see that it compares
 727 very well to the map obtained in experiment 1 (see Fig. 10),
 728 where the motions of the two cameras were constrained by the
 729 stereo rig and a common motion was predicted using odometry.
 730 Fig. 14 bottom shows a detail of the cameras in their final relative
 731 position. We measure an error along the baseline of less than
 732 2 cm. The orientation errors are less than 0.7° .

733 VII. CONCLUSION

734 We showed in this paper that fusing the visual information
 735 with monocular methods while performing multicamera SLAM
 736 provides several advantages: the ability to consider points at in-
 737 finity, desynchronization of the different cameras, the use of any
 738 number of cameras of different types, sensor self-calibration,
 739 and the possibility to conceive decentralized schemes that will
 740 make realistic multirobot monocular SLAM possible. Except for
 741 decentralization, these advantages have been explored with the
 742 inverse depth monocular SLAM algorithm, and applied to two
 743 different problems: stereovision SLAM with an extrinsically
 744 decalibrated stereo rig and cooperative SLAM of two independ-
 745 ently moving cameras.

746 Both demonstrations employed a *master-slave* approach,
 747 which made solving some of the issues of map and image
 748 management easier, and we are now improving on this by im-
 749 plementing a fully symmetrical approach. This approach should
 750 easily permit the extension of the presented applications to cases
 751 with more than two cameras. In parallel to these activities, we
 752 started new work on landmark parametrization to improve EKF
 753 linearity in cases of increasing parallax. Also, as parallax in-
 754 creases, landmarks appearances may change too much as to
 755 guarantee a stable operation with the matching methods pre-
 756 sented here. We believe that wide baseline feature matching
 757 will be the bottleneck of visual SLAM for some time to come.
 758 As for decentralization, we note that it demands a full reformu-
 759 lation of the fusion engines we use in this paper (one central
 760 EKF), for example, via channel filters, and is currently a subject
 761 of intense research at LAAS and other laboratories.

REFERENCES

- [1] J. Solà, A. Monin, M. Devy, and T. Lemaire, “Undelayed initialization in bearing only SLAM,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Edmonton, AB, Canada, Aug. 2–6, 2005, pp. 2499–2504. 763
- [2] A. Chiuso, P. Favaro, H. Jin, and S. Soatto, “Structure from motion causally integrated over time,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 24, no. 4, pp. 523–535, Apr. 2002. 764
- [3] A. J. Davison, “Real-time simultaneous localisation and mapping with a single camera,” in *Proc. Int. Conf. Comput. Vis.*, Nice, France, Oct. 2003, vol. 2, pp. 1403–1410. 765
- [4] J. Civera, A. J. Davison, and J. M. M. Montiel, “Dimensionless monocular SLAM,” in *Proc. Iberian Conf. Pattern Recognit. Image Anal.*, Jun. 2007, pp. 412–419. 766
- [5] J. Civera, A. Davison, and J. Montiel, “Inverse depth parametrization for monocular SLAM,” *IEEE Trans. Robot.*, vol. 24, no. 5, Oct. 2008. 767
- [6] E. Eade and T. Drummond, “Scalable monocular SLAM,” in *Proc. IEEE Int. Conf. Comput. Vis. Pattern Recognit.*, Jun. 17–22, 2006, vol. 1, pp. 469–476. 768
- [7] B. Triggs, P. McLauchlan, R. Hartley, and A. Fitzgibbon, “Bundle adjustment—A modern synthesis,” in *Vision Algorithms: Theory and Practice*, ser. LNCS, W. Triggs, A. Zisserman, and R. Szeliski, Eds. New York: Springer-Verlag, 2000, pp. 298–375. 769
- [8] K. Konolige, “SLAM via variable reduction from constraints maps,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Barcelona, Spain, Apr. 18–22, 2005, pp. 667–672. 770
- [9] J. Folkesson, P. Jensfelt, and H. I. Christensen, “Vision SLAM in the measurement subspace,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Barcelona, Spain, Apr. 18–22, 2005, pp. 30–35. 771
- [10] J. Diebel, K. Reuterswård, S. Thrun, and R. G. J. Davis, “Simultaneous localization and mapping with active stereo vision,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Sendai, Japan, Oct. 2004, vol. 4, pp. 3436–3443. 772
- [11] J. Solà, A. Monin, and M. Devy, “BiCamSLAM: Two times mono is more than stereo,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Rome, Italy, Apr. 2007, pp. 4795–4800. 773
- [12] L. M. Paz, P. Piniés, J. Tardós, and J. Neira, “Large scale 6 DOF SLAM with stereo-in-hand,” *IEEE Trans. Robot.*, vol. 24, no. 5, Oct. 2008. 774
- [13] A. Mallet, S. Lacroix, and L. Gallo, “Position estimation in outdoor environments using pixel tracking and stereovision,” in *Proc. Int. Conf. Robot. Autom.*, San Francisco, CA, 2000, vol. 4, pp. 3519–3524. 775
- [14] K. Konolige, M. Agrawal, and J. Solà, “Large-scale visual odometry for rough terrain,” presented at the Int. Symp. Res. Robot., Hiroshima, Japan, Nov. 2007. 776
- [15] T. D. Barfoot, “Online visual motion estimation using FastSLAM with SIFT features,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Aug. 2–6, 2005, pp. 579–585. 777
- [16] A. I. Comport, E. Malis, and P. Rives, “Accurate quadrifocal tracking for robust 3D visual odometry,” in *Proc. Int. Conf. Robot. Autom.*, Rome, Italy, Apr. 10–14, 2007, pp. 40–45. 778
- [17] E. M. Foxlin, “Generalized architecture for simultaneous localization, auto-calibration, and map-building,” in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, Lausanne, Switzerland, 2002, vol. 1, pp. 527–533. 779
- [18] E. Nettleton, H. Durrant-Whyte, and S. Sukkarieh, “A robust architecture for decentralised data fusion,” presented at the Int. Conf. Adv. Robot., Coimbra, Portugal, 2003. 780
- [19] J. Solà, “Towards visual localization, mapping and moving objects tracking by a mobile robot: a geometric and probabilistic approach” Ph.D. dissertation, Inst. Nat. Polytech. de Toulouse, Toulouse, France, 2007. 781
- [20] J. Civera, A. Davison, and J. Montiel, “Inverse depth to depth conversion for monocular SLAM,” in *Proc. IEEE Int. Conf. Robot. Autom.*, Rome, Italy, Apr. 10–14, 2007, pp. 2778–2783. 782
- [21] A. J. Davison, “Active search for real-time vision,” in *Proc. Int. Conf. Comput. Vis.*, 2005, vol. 1, pp. 66–73. 783
- [22] C. Harris and M. Stephens, “A combined corner and edge detector,” in *Proc. 4th Alvey Vis. Conf.*, Manchester, U.K., 1988, pp. 189–192. 784
- [23] A. J. Davison, I. D. Reid, N. D. Molton, and O. Stasse, “MonoSLAM: Real-time single camera SLAM,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 29, no. 6, pp. 1052–1067, Jun. 2007. 785
- [24] N. Molton, A. J. Davison, and I. Reid, “Locally planar patch features for real-time structure from motion,” presented at the Brit. Mach. Vis. Conf., Kingston, U.K., 2004. 786
- [25] K. Strobl, W. Sepp, S. Fuchs, C. Paredes, and K. Arbter. (2006). Camera calibration toolbox for Matlab. Inst. Robot. Mechatronics, Wessling, Germany, Tech. Rep. [Online]. Available: http://www.vision.caltech.edu/bouguetj/calib_doc/index.html 787

836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851

Joan Solà was born in Barcelona, Spain, in 1969. He received the B.Sc. degree in telecommunications and electronic engineering from the Universitat Politècnica de Catalunya, Barcelona, in 1995, the M.Sc. degree in control systems from the École Doctorale Systèmes, Toulouse, France, in 2003, and the Ph.D. degree in control systems from the Institut National Polytechnique de Toulouse in 2007, where he was hosted by the Laboratoire d'Analyse et d'Architecture des System (LAAS), Centre National de la Recherche Scientifique (CNRS).

He was a Postdoctoral Fellow at SRI International, Menlo Park, CA. He is currently at LAAS-CNRS, where he is engaged in research on visual localization and mapping. His current research interests include estimation and data fusion applied to off-road navigation, mainly using vision.



Michel Devy received the degree in computer science engineering from the Ecole Nationale Supérieure d'Ingénieurs Electriciens de Grenoble, Grenoble, France, in 1976 and the Ph.D. degree from the Laboratoire Automatique et d'Analyse des Systèmes, Centre Nationale de la Recherche Scientifique (LAAS-CNRS), Toulouse, France, in 1980.

Since 1980, he has been with the Department of Robotics and Artificial Intelligence, LAAS-CNRS, where he is the Research Director and the Head of the Research Group Robotics, Action, and Perception. His current research interests include computer vision for automation and robotics applications. He has also been involved in numerous national and international projects concerning, about field and service robots, 3-D vision for intelligent vehicles, 3-D metrology, and others. He has authored or coauthored about 150 scientific communications.

868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867

André Monin was born in Le Creusot, France, in 1958. He received the Graduate degree from the Ecole Nationale Supérieure d'Ingénieurs Electriciens de Grenoble, Grenoble, France, in 1980, the Ph.D. degree in nonlinear systems representation from the University Paul Sabatier, Toulouse, France, in 1987, and the Habilitation pour Diriger des Recherches degree from the University Paul Sabatier in 2002.

From 1981 to 1983, he was a Teaching Assistant with the Ecole Normale Supérieure de Marrakech, Marrakech, Morocco. Since 1985, he has been with the Laboratoire Automatique et d'Analyse des Systèmes, Centre Nationale de la Recherche Scientifique (LAAS-CNRS), Toulouse, as the "Chargé de Recherche." His current research interests include the areas of nonlinear filtering, systems realization, and identification.



Teresa Vidal-Calleja received the B.Sc. degree in mechanical engineering from the Universidad Nacional Autónoma de México, México City, México, in 2000, the M.Sc. degree in mechatronics from CINVESTAV-IPN, México City, in 2002, and the Ph.D. degree in robotics, automatic control, and computer vision from the Universitat Politècnica de Catalunya, Barcelona, Spain, in 2007.

She was a Visiting Research Student with the University of Oxford's Robotics Research Group, the Australian Centre for Field Robotics, and the University of Sydney. She is currently a Postdoctoral Fellow with the Robotics and Artificial Intelligence Group, Laboratoire Automatique et d'Analyse des Systèmes, Centre Nationale de la Recherche Scientifique (LAAS-CNRS), Toulouse, France. Her current research interests include autonomous vehicles, perception, control, and cooperation.

885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901

Q2

Q3

QUERIES

- 903 Q1: Author: Please reframe the references to colors in the caption of Figs. 5, 8, 9, 11, 12, and 14, if the artwork is not being
904 produced in color
- 905 Q2. Author: Please check if the details of the academic degrees received by A. Monin are OK as edited.
- 906 Q3. Author: Please provide the title of first degree. Also, provide the subject (physics, mathematics, electrical engineering, etc.)
907 in which M. Dey received the Ph. D. degree.

IEEE
PROOF