

# Pose-graph underwater Simultaneous Localization And Mapping for autonomous monitoring by means of optical and acoustic sensors

Alessandro Bucci<sup>1</sup>, Alessandro Ridolfi<sup>1</sup>, and Benedetto Allotta<sup>1</sup>

<sup>1</sup>Universita degli Studi di Firenze Dipartimento di Ingegneria Industriale

August 26, 2023

## Abstract

Modern mobile robots require precise and robust localization and navigation systems to achieve mission tasks correctly. In particular, in the underwater environment, where Global Navigation Satellite Systems (GNSSs) cannot be exploited, the development of localization and navigation strategies becomes more challenging. Maximum A Posteriori (MAP) strategies have been analyzed and tested to increase navigation accuracy and take into account the entire history of the system state. In particular, a sensor fusion algorithm relying on a MAP technique for Simultaneous Localization and Mapping (SLAM) has been developed to fuse information coming from a monocular camera and a Doppler Velocity Log (DVL) and to consider the landmark points in the navigation framework. The proposed approach can guarantee to simultaneously locate the vehicle, thanks to the onboard sensors, and map the surrounding environment with the information extracted from the images acquired by a bottom-looking optical camera. Optical sensors can provide constraints between the vehicle poses and the landmarks belonging to the observed scene. The DVL measurements have been employed to solve the unknown scale factor and to guarantee the correct vehicle localization even in absence of visual features. After validating the solution through realistic simulations, an experimental campaign at sea was conducted in Stromboli Island (Messina), Italy. In conclusion, an algorithm, which works with the Poisson surface reconstruction method to obtain a smooth seabed surface, for mesh creation has been developed.

# Pose-graph underwater Simultaneous Localization And Mapping for autonomous monitoring by means of optical and acoustic sensors

---

**Alessandro Bucci\***

Department of Industrial Engineering  
University of Florence  
via di Santa Marta 3  
50139, Florence, Italy  
Interuniversity Center of  
Integrated Systems for the  
Marine Environment (ISME)

**Alessandro Ridolfi**

Department of Industrial Engineering  
University of Florence  
via di Santa Marta 3  
50139, Florence, Italy  
Interuniversity Center of  
Integrated Systems for the  
Marine Environment (ISME)

**Benedetto Allotta**

Department of Industrial Engineering  
University of Florence  
via di Santa Marta 3  
50139, Florence, Italy  
Interuniversity Center of  
Integrated Systems for the  
Marine Environment (ISME)

## Abstract

Modern mobile robots require precise and robust localization and navigation systems to achieve mission tasks correctly. In particular, in the underwater environment, where Global Navigation Satellite Systems (GNSSs) cannot be exploited, the development of localization and navigation strategies becomes more challenging. Maximum A Posteriori (MAP) strategies have been analyzed and tested to increase navigation accuracy and take into account the entire history of the system state. In particular, a sensor fusion algorithm relying on a MAP technique for Simultaneous Localization and Mapping (SLAM) has been developed to fuse information coming from a monocular camera and a Doppler Velocity Log (DVL) and to consider the landmark points in the navigation framework. The proposed approach can guarantee to simultaneously locate the vehicle, thanks to the onboard sensors, and map the surrounding environment with the information extracted from the images acquired by a bottom-looking optical camera. Optical sensors can provide constraints between the vehicle poses and the landmarks belonging to the observed scene. The DVL measurements have been employed to solve the unknown scale factor and to guarantee the correct vehicle localization even in absence of visual features. After validating the solution through realistic simulations, an experimental campaign at sea was conducted in Stromboli Island (Messina), Italy. In conclusion, an algorithm, which works with the Poisson surface reconstruction method to obtain a smooth seabed surface, for mesh creation has been developed.

---

\*Corresponding author, email address: [alessandro.bucci@unifi.it](mailto:alessandro.bucci@unifi.it)

# 1 Introduction

From geology to exploration and surveillance of archaeological sites and from Oil and Gas industry to reconnaissance for military purposes, exploring and understanding seas and oceans is a matter of primary importance. Considering their human hostile nature, since the 1960s, seas and oceans have been explored with the aid of robots. The first Unmanned Underwater Vehicles (UUVs) were teleoperated ones and are referred in the technical literature as Remotely Operated Vehicles (ROVs). A cable, usually called umbilical cable, acts as a constant connection providing power and communications, and specialised operators are thus able to control the vehicle using the feedback forwarded by the on-board sensors. In the last decades Autonomous Underwater Vehicles (AUVs), which are completely autonomous, have gained interest with respect to ROVs. Indeed, such vehicles do not require human intervention (except for deployment and recovery), are usually equipped with electric batteries, and possess dedicated systems used to control their motion. Since the demanded tasks of underwater vehicles have become more and more challenging (Prats et al., 2012), (Ferri et al., 2017), researchers and scientists are following the tide of change and are pushing the boundaries of AUVs capabilities by integrating cutting-edge technologies. Indeed, autonomous inspection (Cashmore et al., 2014), and intervention (Youakim et al., 2020) strategies for underwater installations, exploration planning solutions (Vidal et al., 2020), and autonomous coverage approaches (Paull et al., 2012), have become essential tools to execute demanding and hazardous subsea operations.

One of the most significant and complex tasks in autonomous underwater exploration is to retrieve the vehicle's pose within the surrounding environment, making use of precise and reliable navigation and localization systems, which are necessary regardless of the kind of mission or task the underwater vehicle is required to perform. In addition to this, perceptual devices (such as optical cameras and acoustic devices) able to sense the surrounding environment have been earning attention throughout the last decades to acquire data for monitoring and inspection purposes. The use of optical and acoustic equipment to aid navigation has emerged as a relevant alternative or support to traditional navigation sensors.

Several algorithms have been developed throughout the years to increase the navigation and localization capabilities of the AUVs relying on Bayesian estimators, such as Kalman filtering and Maximum A Posteriori (MAP) estimators. Both Extended Kalman Filter (EKF) (Dissanayake et al., 2001) and least squares optimization (Dellaert and Kaess, 2006) have been used extensively in Simultaneous Localization And Mapping (SLAM) research in the past (Zhang et al., 2018). Earlier SLAM research has used EKF algorithms where the state vector contained the latest robot pose and the positions of the observed features. However, it has been shown that EKF-SLAM could result in inconsistent estimate (Julier and Uhlmann, 2001), (Castellanos et al., 2004), as the estimated covariance from the algorithm can violate the theoretical achievable lower bounds (Dissanayake et al., 2001), (Huang and Dissanayake, 2007). On the contrary, optimization based SLAM uses a state vector containing all the robot poses and all the features observed. Considering that relinearization is performed during each iteration step, there is no inconsistency issue in optimization based SLAM and thus the quality of the estimate is higher than that of EKF-SLAM.

Consequently, to overcome the limitations introduced by the Kalman filter strategies, which condense all the history into the last estimation, a sensor fusion MAP algorithm has been developed for underwater navigation in the context of this work. Due to the complexity of retrieving navigation information in the underwater environment, a sensor fusion approach has been used. The performance and robustness of the visual SLAM algorithm heavily rely on the quality of the images and salient features. Consequently, the visual SLAM system has been fused with other sensing algorithms, such as the Doppler Velocity Log (DVL). As shown previously, very few works still exist on underwater SLAM fusing data from a monocular camera and a DVL. Despite that, fusing an optical and an acoustic sensor in a MAP-based framework can take advantages from both sensors, which have an excellent complement to each other. This developed solution can be employed to locate the vehicle and map the seabed at the same time in a unified framework. Thus, an underwater visual acoustic SLAM strategy which integrates DVL with a visual SLAM system has been developed to perform accurate navigation and mapping tasks at the same time. Particular attention has been focused on the design of scale factor ambiguity resolution and extrinsic calibration optimization procedure and on implementing a reset procedure to reduce the computational burden. Furthermore, the proposed strategy has been tested with both simulated and experimental data to evaluate the navigation performance and has been compared with an Unscented Kalman Filter (UKF)-based algorithm, whose performance has

53 been accurately discussed in authors' previous works (Bucci et al., 2023), (Bucci et al., 2021).  
54 The paper is organized as follows: state-of-the-art in SLAM strategies are detailed in Section 2, whereas  
55 Section 3 is dedicated to introduce the MAP estimation approach. Section 4 outlines the development of the  
56 factor graph framework, whereas some improvements and peculiarity of the proposed SLAM strategy are  
57 reported in Section 5. While navigation results obtained from simulated environment and from an experi-  
58 mental campaign are reported respectively in Section 6 and 7, an analysis of the mapping capabilities are  
59 depicted in Section 8. Finally, Section 9 draws conclusions.

## 60 2 Related works

61 Many estimation problems in robotics have an underlying optimization problem (Dellaert, 2021). In most  
62 of these optimization problems, the objective to be maximized or minimized is composed of many differ-  
63 ent factors (e.g., a Global Navigation Satellite System (GNSS) measurement is applied to the pose of the  
64 vehicle at a particular time and can be referred as an unary factor, an Inertial Measurement Unit (IMU)  
65 measurement can be related to two vehicle states at adjacent times and can represent an odometry factor).  
66 The use of factorial graphs in the design of algorithms for robotic applications has three main advantages.  
67 First, since many optimization problems in robotics have the property of locality, factorial graphs can model  
68 a wide variety of problems in all robotics domains, such as tracking, navigation, and mapping. Secondly,  
69 by clearly exposing the structure of the problem, reflection on factorial graphs offers many opportunities to  
70 improve the performance of key algorithms. Many classical algorithms can be viewed as the application of  
71 the elimination algorithm to a particular type of factorial graph. Still, this algorithm is only optimal for a  
72 small class of problems. In many applications, knowledge of the specific structure of the problem domain  
73 can improve the execution time of inference by orders of magnitude. Similarly, well-known algorithmic ideas  
74 from linear algebra can be generalized to factorial graphs, leading, for example, to incremental inference  
75 algorithms. Thirdly, apart from performance considerations, factorial graphs are useful when designing and  
76 thinking about how to model a problem, providing a common language to express ideas to collaborators and  
77 users of a particular algorithm. After working with factor graphs for a while, one begins to identify factor  
78 types as a particularly useful design unit. A factor type specifies how many variables a factor is related to  
79 and the semantics associated with the function to be calculated.

80 MAP estimation has recently become the standard approach for modern SLAM strategies (Cadena et al.,  
81 2016). Indeed, while fixed-lag smoothers and filtering solutions restrict the inference within a window of the  
82 latest states or to the latest state, respectively, MAP strategies estimate the entire history of the system  
83 by solving a non-linear optimization problem. Both fixed-lag smoothers and filters marginalize older states,  
84 collapsing the corresponding information (usually) in a Gaussian prior. This approach can lead to reduced  
85 robustness against outlier data (Forster et al., 2016). Since MAP strategies can quickly lead to an unsuitable  
86 approach for real-time applications, the development of incremental smoothing techniques has arisen as the  
87 state-of-the-art approach. Such techniques can reuse previously calculated quantities when new measure-  
88 ments or variables are added (Kaess et al., 2008), (Kaess et al., 2012). In particular, in (Kaess et al., 2012)  
89 a Bayes tree data structure is employed to perform incremental optimization on the factor graph. Also, the  
90 adopted solution possesses the ability to identify and update only a small subset of variables by accurately  
91 selecting the ones affected by the new measurement. A complete review can be found in (Grisetti et al.,  
92 2020) and the references therein.

93 Considering the underwater domain, two works have been taken as inspiration for the development of the  
94 factor graph employed in the proposed SLAM strategy. (Westman and Kaess, 2019) proposes an algorithm  
95 to generate pose-to-pose constraints for pairs of SONAR images and to fuse these resulting pose constraints  
96 with the vehicle odometry in a pose graph optimization framework. In (Franchi et al., 2021) Ultra-Short  
97 BaseLine (USBL) measurements are exploited as observations within the on-board navigation filter, where  
98 the localization task is solved as a MAP estimation problem. Both these solutions rely on Incremental  
99 Smoothing and Mapping 2 (iSAM2), which is the last evolution of the incremental smoothing and mapping  
100 solution developed in Georgia Tech Smoothing And Mapping (GTSAM). Furthermore, other graph-based  
101 SLAM strategies have been proposed to fuse the data obtained by the navigation sensors and the perception  
102 sensors, both acoustic and optical. In (Fallon et al., 2013) this approach is used in an AUV for mine counter

103 measurement and localization. While the graph is initialized by pose node from a Global Positioning System  
 104 (GPS), a non-linear least square optimization is performed with the DVL and IMU-based Dead Reckoning  
 105 (DR) estimations and the SONAR images. In (Huang and Kaess, 2015) an acoustic structure from motion  
 106 algorithm for recovering 3D scene structure from multiple 2D SONAR images while at the same time local-  
 107 izing the SONAR is presented.

108 Turning to visual SLAM, ORB-SLAM (Mur-Artal et al., 2015) is one of the most complete and simple  
 109 algorithms, and the whole system is calculated around Oriented FAST and Rotated BRIEF (ORB) fea-  
 110 ture points, with features such as rotational scale invariance and fast detection. ORB-SLAM2 (Mur-Artal  
 111 and Tardós, 2017) is upgraded from ORB-SLAM, supporting monocular, binocular, and RGB-D modes,  
 112 and has good adaptability. Finally, the latest ORB-SLAM3 (Campos et al., 2021) algorithm fuses opti-  
 113 cal images with inertial sensors. The excellent characteristics of the ORB-SLAM2 algorithm, which can  
 114 achieve centimeter-level precision on the ground, represent an incentive for its application in underwater  
 115 environments. Consequently, the visual part of the developed SLAM algorithm takes inspiration from the  
 116 ORB-SLAM2 framework. Referring to the vision-based SLAM algorithm for underwater navigation and  
 117 mapping, (Hong and Kim, 2020) addresses a visual mapping method for precise camera trajectory estima-  
 118 tion and 3D reconstruction of underwater ship hull surface using a monocular camera as the primary sensor.  
 119 (Du et al., 2017) proposes an underwater visual SLAM system using a stereo camera, which has been tested  
 120 in a circular pool.

121 Finally, an acoustic-visual-inertial SLAM strategy has been proposed in (Rahman et al., 2018) and (Rah-  
 122 man et al., 2018). Data coming from a mechanical scanning SONAR, a stereo camera, and proprioceptive  
 123 inertial sensors are fused in a tightly coupled non-linear optimization to estimate the vehicle trajectory and  
 124 reconstruct the surrounding environment. There are few works where the DVL measurements are fused with  
 125 other perception sensors in a SLAM strategies. In (Ozog and Eustice, 2013) a SLAM method, which uses a  
 126 very sparse point cloud derived from a DVL to add constraints to a piecewise-planar framework, is proposed.  
 127 A camera is also employed to bound drifts of odometry fused by a DVL, IMU and pressure Depth Sensor  
 128 (DS) (Kim and Eustice, 2013). Fiducial markers are also integrated into a visual SLAM framework with  
 129 DVL, IMU, and DS in (Westman and Kaess, 2018).

### 130 3 Maximum A Posteriori estimation

131 A navigation and mapping problem is a problem where the unknown state variables  $X = \{x_1, x_2, \dots, x_M\}$   
 132 constituted of poses and landmarks has to be determined given the measurements  $Z = \{z_1, z_2, \dots, z_N\}$ . The  
 133 MAP estimator maximizes the posterior density  $p(X|Z)$  of the states  $X$  given the measurements  $Z$ :

$$134 \quad X^{MAP} = \underset{X}{\operatorname{argmax}} p(X)l(Z|X) = p(X) \prod_{i=1}^N l(z_i|X), \quad (1)$$

135 where  $l(z_i|X)$  is the likelihood distribution and an additive Gaussian noise is assumed in all measurement  
 136 models, as reported in Eq. 2.

$$137 \quad p(z_i|X) = \mathcal{N}(h_i(X), \Sigma_i) \propto \exp\left(-\frac{1}{2}\|h_i(X) - z_i\|_{\Sigma_i}^2\right) \quad (2)$$

138 where  $h_i(X)$  is the measurement function, which maps the state estimate  $X$  into a predicted value  $\hat{z}_i$  of  
 139 the measurement  $z_i$  and  $\Sigma_i$  is the covariance matrix, which summarizes the uncertainty of the measurement  
 140 model. By applying the monotonic logarithmic function and the Gaussian model previously introduced, the  
 141 optimization problem can be simplified into a nonlinear least square problem:

$$142 \quad X^{MAP} = \underset{X}{\operatorname{argmin}} \sum_{i=1}^N \|h_i(X) - z_i\|_{\Sigma_i}^2 \quad (3)$$

143 where

$$144 \quad \|h_i(X) - z_i\|_{\Sigma_i}^2 = (h_i(X) - z_i)^\top \Sigma_i^{-1} (h_i(X) - z_i) \quad (4)$$

145 is the Mahalanobis distance.

146 The nonlinear problem can be solved through standard methods, such as the Gauss-Newton or the Levenberg-  
 147 Marquardt algorithms, which iteratively converge to the solution by solving the linear approximation of the  
 148 nonlinear system. More information can be found in (Grisetti et al., 2020), (Dellaert and Kaess, 2017).

## 149 4 Factor graph framework development

150 The mathematical modeling of the factors used to represent the measurement constraints to solve the au-  
 151 tonomous navigation and mapping problem is presented. Inspired by (Westman and Kaess, 2018), (Westman  
 152 and Kaess, 2020), the factors described below have been employed, where it is necessary to consider that the  
 153 information included in some factors can be derived from measurements not coming from a single sensor.  
 154 The state of the system at instant  $i$  is defined as a complete pose belonging to SE(3), which can be expressed  
 155 mathematically as:

$$156 \quad T_{x_i} = \begin{bmatrix} R_i & \mathbf{t}_i \\ \mathbf{0}^{1 \times 3} & 1 \end{bmatrix} \quad (5)$$

157 where  $R_i \in \text{SO}(3)$  is the rotation matrix and  $\mathbf{t}_i \in \mathbb{R}^3$  represents the translation vector. Defining the set of  
 158 poses at time  $k$  with  $\mathcal{X}_k$ , such that  $\mathcal{X}_k = \{T_{x_i}\}_{i=0,1,\dots,k}$ , it is possible to define the optimization problem  
 159 and, in particular, Eq. 4 on the smooth manifold SE(3). Considering a transformation from the state  $x_i$  to  
 160 the state  $x_j$  constrained with an odometry measurement  $z_{i,j}$  with covariance  $\Sigma_{i,j}$ , Eq. 4 becomes:

$$161 \quad \|f_{ij}(x_i, x_j) \ominus z_{i,j}\|_{\Sigma_{i,j}}^2 = \|\log(T_{z_{i,j}}^{-1} T_{x_i}^{-1} T_{x_j})\|_{\Sigma_{i,j}}^2 \quad (6)$$

162 The symbol  $\ominus$  encodes the logarithmic map from the manifold to an element of the SE(3) Lie algebra, where  
 163  $f_{ij}(\cdot)$  represents the measurement function applied to the poses  $T_{x_i}$  and  $T_{x_j}$ . For ease of explanation  $T_{x_i}$   
 164 can be represented with the vector  $[X_{x_i} \ Y_{x_i} \ Z_{x_i} \ \phi_{x_i} \ \theta_{x_i} \ \psi_{x_i}] \in \mathbb{R}^6$  and the measurement function  
 165 becomes

$$166 \quad f_{ij}(x_i, x_j) = [X_{x_{i,j}} \ Y_{x_{i,j}} \ Z_{x_{i,j}} \ \phi_{x_{i,j}} \ \theta_{x_{i,j}} \ \psi_{x_{i,j}}]^\top \quad (7)$$

167 In contrast, for a measurement  $z_i$  that indicates a local information on the state  $x_i$  with covariance  $\Sigma_i$ , Eq.  
 168 4 is

$$169 \quad \|f_i(x_i) \ominus z_i\|_{\Sigma_i}^2 = \|\log(T_{z_i}^{-1} T_{x_i})\|_{\Sigma_i}^2 \quad (8)$$

170 where the measurement function  $f_i(\cdot)$  applied to the pose  $T_{x_i}$  can be defined as:

$$171 \quad f_i(x_i) = [X_{x_i} \ Y_{x_i} \ Z_{x_i} \ \phi_{x_i} \ \theta_{x_i} \ \psi_{x_i}]^\top \quad (9)$$

172 The information from the available onboard sensors has been encoded as measurement factors to constrain  
 173 the optimization, whose solution represents the MAP estimate. Inspired by (Westman and Kaess, 2019), the  
 174 following factors have been included:

- 175 • a relative 4D pose-to-pose constraint on  $x$ ,  $y$ , and  $z$  translation and yaw rotation, thanks to the  
 176 measurements coming from the DVL and the yaw estimated by the attitude estimator;
- 177 • a unary 2D constraint on pitch and roll rotations, obtained from the attitude estimation filter;
- 178 • a unary 1D constraint on  $z$  translation thanks to the DS measurements;
- 179 • a unary constraint on  $x$  and  $y$  translation exploiting GNSS observations;
- 180 • a relative 6D pose-to-pose constraint on  $x$ ,  $y$ , and  $z$  translation and roll, pitch, and yaw rotation,  
 181 thanks to the relative pose estimated through the monocular camera and properly scaled;
- 182 • a camera-based landmark constraint on the vehicle pose and the landmark position for each feature  
 183 seen with the monocular camera over multiple images.

184 The implemented approach adds a new state only when at least one observation from GNSS, DVL, DS, or  
 185 when the visibility is acceptable, the camera is available. The link between adjacent nodes is maintained  
 186 by collapsing the relative motion XYZ-Y in a single compound constraint, where simple DR is performed  
 187 between the two consecutive nodes with the last acquired DVL measurements. The pose  $T_{x_i}$  can be repre-  
 188 sented with a vector  $[X_{x_i} \ Y_{x_i} \ Z_{x_i} \ \phi_{x_i} \ \theta_{x_i} \ \psi_{x_i}] \in \mathbb{R}^6$  that encodes the state at the generic instant.  
 189 Mathematically, at time  $k$ , the optimization problem can be written as

$$190 \quad \mathcal{X}_k^* = \underset{X}{\operatorname{argmax}} \sum_{i=1}^{k-1} \left( \|m_{XYZ-Y}(x_{i-1}, x_i) \ominus o_{i-1,i}\|_{\Sigma_{o_{i-1,i}}}^2 + \|m_{RP}(x_i) \ominus r_i\|_{\Sigma_{r_i}}^2 \right) + \quad (10)$$

$$+ \sum_{i \in \mathcal{Z}} \|m_Z(x_i) - z_i\|_{\Sigma_{z_i}}^2 +$$

$$+ \sum_{i \in \mathcal{G}} \|m_{XY}(x_i) - \mathbf{g}_i\|_{\Sigma_{\mathbf{g}_i}}^2 +$$

$$+ \sum_{i,j \in \mathcal{C}} \|m_{XYZ-RPY}(x_i, x_j) \ominus p_{i,j}\|_{\Sigma_{p_{i,j}}}^2 +$$

$$+ \sum_{j \in \mathcal{LM}, i \in \mathcal{C}} \rho \left( \|\mathbf{p}_{ij} - \pi_i(T_{x_i} \mathbf{P}_j)\|_{\Sigma_{lm_i}}^2 \right) +$$

$$+ \|T_{x_0} \ominus T_{x_{prior}}\|_{\Sigma_{lm_i}}^2$$

191  $\{m_{XYZ-Y}(\cdot), o_{i-1,i}, \Sigma_{o_{i-1,i}}\}$ ,  $\{m_{RP}(\cdot), r_i, \Sigma_{r_i}\}$ ,  $\{m_Z(\cdot), z_i, \Sigma_{z_i}\}$ ,  $\{m_{XY}(\cdot), \mathbf{g}_i, \Sigma_{\mathbf{g}_i}\}$ ,  
 192  $\{m_{XYZ-RPY}(\cdot), p_{i,j}, \Sigma_{p_{i,j}}\}$  are the measurement functions, the measured values and covariances as-  
 193 sociated to the previously introduced factors. In particular,  $o_{i-1,i}$ ,  $r_i$  and  $p_{i,j}$  represent, on SE(3), the  
 194 observation for the XYZ-Y part, the RP part and the camera-based XYZ-RPY part, respectively,  $z_i \in \mathbb{R}$   
 195 is the depth measurement,  $\mathbf{g}_i \in \mathbb{R}^2$  is the GNSS measurement. The measurement functions are:

$$m_{XYZ-Y}(x_{i-1}, x_i) = [X_{x_{i-1,i}} \ Y_{x_{i-1,i}} \ Z_{x_{i-1,i}} \ \psi_{x_{i-1,i}}]^\top$$

$$196 \quad m_{RP}(x_i) = [\phi_{x_i} \ \theta_{x_i}]^\top$$

$$m_Z(x_i) = [Z_{x_i}]$$

$$m_{XY}(x_i) = [X_{x_i} \ Y_{x_i}]^\top$$

$$m_{XYZ-RPY}(x_i, x_j) = [X_{x_{i,j}} \ Y_{x_{i,j}} \ Z_{x_{i,j}} \ \phi_{x_{i,j}} \ \theta_{x_{i,j}} \ \psi_{x_{i,j}}]^\top$$

197 Thanks to the features extracted from optical images and matched through multiple keyframes, it is possible  
 198 to optimize map point locations  $\mathbf{P}_j \in \mathbb{R}^3$  and keyframe poses  $T_{x_i} \in \text{SE}(3)$  minimizing the reprojection error  
 199 with respect to the matched keypoints  $\mathbf{p}_{ij} \in \mathbb{R}^2$ . The error term for the observation of a map point  $j$  in a  
 200 keyframe  $i$  is

$$201 \quad \mathbf{e}_{ij} = \mathbf{p}_{ij} - \pi_i(T_{x_i} \mathbf{P}_j) \quad (12)$$

202 where  $\pi_i(\cdot)$  is the projection function:

$$203 \quad \pi_i(T_{x_i} \mathbf{P}_j) = \begin{bmatrix} f_x \frac{x_{ij}}{z_{ij}} + c_x \\ f_y \frac{y_{ij}}{z_{ij}} + c_y \end{bmatrix} \quad (13)$$

204 where  $(f_x, f_y)$  and  $(c_x, c_y)$  are respectively the focal length and the principal point of the camera and  
 205  $[x_{ij} \ y_{ij} \ z_{ij}]^\top$  are the coordinates of the point. The cost function to be minimized can be defined as:

$$206 \quad f_{LM}(x_i) = \rho \left( \|\mathbf{p}_{ij} - \pi_i(T_{x_i} \mathbf{P}_j)\|_{\Sigma_{lm_i}}^2 \right) \quad (14)$$

207 where  $\rho(\cdot)$  is the Huber robust cost function and  $\Sigma_{lm_i}$  is the covariance matrix associated to the scale at  
 208 which the keypoint  $i$  was detected. While  $\mathcal{Z}$ ,  $\mathcal{G}$  and  $\mathcal{C}$  are the set of pose nodes for which DS, GNSS and  
 209 camera measurements respectively occur,  $\mathcal{LM}$  is the set of landmark nodes.  $T_{x_{prior}}$  is the prior constraint  
 210 on the first pose, which is necessary to anchor the state evolution to a global coordinate frame (Fig. 1).

211 In terms of implementation, the GTSAM library (Dellaert, 2012) has been used as the back-end for the  
 212 localization solution. Further information can be found in (Kaess et al., 2008), (Kaess et al., 2012). iSAM2,  
 213 which is the latest evolution of the incremental smoothing and mapping solution developed in GTSAM,  
 214 allows only the typical small subset of variables affected by a new measurement, i.e., the measurement func-  
 215 tion and associated covariances, to be identified and updated, thus limiting the computational load of the

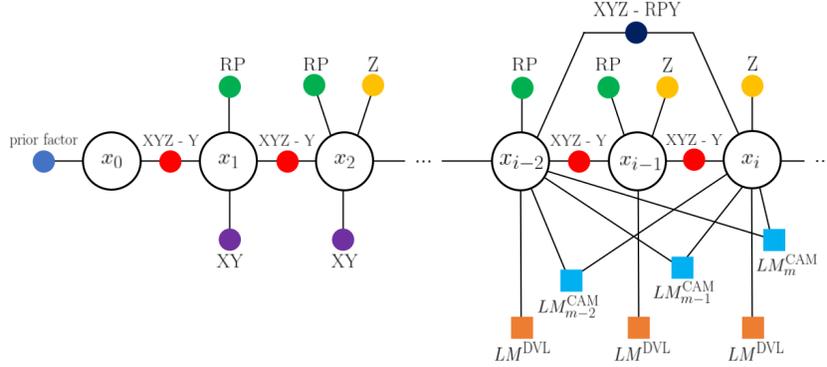


Figure 1: Example of the factor graph at the iteration  $i$  constrained with vision-based landmarks and all the onboard sensors.

216 estimation, offering a trade-off between accuracy and efficiency. Several issues affect the vision in underwa-  
 217 ter environments, which can negatively influence the employment of visual SLAM algorithms. Specifically,  
 218 while scattering reduces light intensity causing a loss of contrast and haze in underwater images, light ab-  
 219 sorption leads to a decrease in the color quality of underwater images. Light attenuation in water introduces  
 220 degradation in underwater images, such as poor colors, decreased contrast, haziness, and blurring, making  
 221 them hardly usable for the filter. Thus it is necessary to guarantee that the visual part of the navigation  
 222 framework, which is dependent on uncontrollable external conditions, can be correctly inserted or removed  
 223 from the factor graph. Only when the visual SLAM algorithm is correctly initialized and the current scale  
 224 factor is computed, it is possible to insert keyframe poses and map point locations in the factor graph. Edges  
 225 computed thanks to other onboard sensors, which do not suffer from visibility limitations, are inserted in  
 226 the whole factor graph during the entire mission. The developed system, through the map points obtained  
 227 from the vision system and the DVL beams, can build a map of the surrounding environment independently  
 228 from the visibility. Indeed, when the reduced visibility impedes the usage of the visual-based features as  
 229 map points, DVL-based beams can be employed to build an approximated map of the sea bottom. The  
 230 quality and resolution of the produced map depend on the availability of the visual landmarks. Still, thanks  
 231 to the DVL measurements, the reconstruction can be performed for the whole mission. Considering that,  
 232 when an AUV accomplishes an underwater mission, the sea bottom texture can change very fast, and its  
 233 depth can increase rapidly, the possibility to guarantee reconstruction of the surrounding environment, even  
 234 approximated, represents a helpful advantage. Obviously, it is necessary to highlight that the DVL beams  
 235 cannot be employed as landmark nodes in the factor graph. Still, they can only be added to the map utilizing  
 236 the sensor geometrical model. It is necessary to highlight that underwater SLAM fusing camera, and DVL  
 237 sensors can increase the localization accuracy and robustness thanks to the excellent complement between  
 238 these two sensors: DVL provides reliable motion estimates for underwater visual SLAM, extending SLAM's  
 239 robustness and operation even without visual features, and vision, when applicable, helps the estimation  
 240 process by introducing visual landmarks which increase the constraints on the vehicle position.  
 241 Turning to the specific strategies for DVL and camera-based factor graph constraint computation, the fol-  
 242 lowing approaches have been employed. It is necessary to notice that while the primary application field  
 243 of a DVL is vehicle navigation through a DR strategy that computes the AUV position by integrating the  
 244 measured linear velocity, the DVL has four acoustic beams, each pointing in a different direction, which can  
 245 be employed to acquire the 3D location of 4 points of the sea bottom during each speed measurement. The  
 246 points located thanks to the DVL beams cannot be employed as additional constraints in the navigation pose  
 247 graph because they do not link any node of the graph. Still, they can easily be used to increase the number  
 248 of points in the estimated map of the sea bottom. Indeed, by knowing the vehicle's actual position from the  
 249 navigation algorithm, the location of the four beams can be converted from the DVL frame to the North,  
 250 East, Down (NED) reference system. The visual SLAM algorithm employed in the developed navigation  
 251 framework is a feature-based monocular SLAM system that operates to estimate the camera trajectory and  
 252 an environment map. The basic idea of the SLAM system introduced in the navigation filter takes inspira-  
 253 tion from the algorithms proposed in (Mur-Artal et al., 2015), (Mur-Artal and Tardós, 2017). Furthermore,

254 following the results reported in (Zacchini et al., 2019), (Bucci et al., 2022), where accurate comparisons  
 255 between several feature detectors are explained, ORB feature detector has been chosen as the preferable  
 256 solution instead of Scale Invariant Feature Transform (SIFT), Speeded Up Robust Features (SURF) and  
 257 Accelerated-KAZE (AKAZE). Considering that a monocular camera is employed, a scale factor ambiguity  
 258 to be solved features the visual-based estimate.

## 259 5 Factor graph framework improvements

### 260 5.1 Scale factor ambiguity resolution

261 This procedure, which is executed every time the visual SLAM algorithm is correctly initialized, has two  
 262 main purposes, the scale factor ambiguity resolution and accurate compensation of the fixed roto-translation  
 263 between the camera and the body frames. This transformation is represented as a similarity transforma-  
 264 tion composed of a scale factor  $s$ , a translation vector  $\mathbf{t}_{c,b} = [t_{c,b}^x \ t_{c,b}^y \ t_{c,b}^z]^\top$  and a rotation matrix  
 265  $R_c^b = R_z(\psi_c^b)R_y(\theta_c^b)R_x(\phi_c^b)$ . It is based on comparing the trajectories estimated through the DVL and the  
 266 other inertial sensors and the camera. It is necessary to notice that until the scale factor has not been esti-  
 267 mated, the measurements obtained thanks to the visual SLAM algorithm are not inserted in the whole factor  
 268 graph. Considering this algorithm’s two purposes and that, usually, underwater vehicles for survey missions  
 269 execute planar trajectories at constant depth, the problem has been solved with a two-step algorithm. In  
 270 particular, while the first part of the algorithm determines a closed-form solution for the  $x$  and  $y$  directions,  
 271 yaw rotation, and the scale factor, the second part optimizes the whole scaled roto-translation with an itera-  
 272 tive algorithm. This framework has been adopted due to the limitations introduced by the particular motion  
 273 executed by the AUV. Indeed, on the one hand, the optimal closed-form solution estimated with 3D points  
 274 that almost lie on a plane cannot correctly estimate the roll and pitch angles of the rigid transformation  
 275 between the two considered reference frames. On the other hand, the iterative algorithm locally converges  
 276 and requires an initial guess in the neighborhood of the exact solution, which can be measured directly on  
 277 the vehicle or evaluated through the closed-form solution.

278 The two steps of the algorithm are described in detail. Firstly, the closed-form solution is found by comput-  
 279 ing the trajectory alignment transformation with translational component on the  $xy$ -plane of the trajectory  
 280 estimated with the DVL and the camera and with rotational component computed with respect to the  
 281 perpendicular axis to this plane. Given the DVL-based positions  $\{\mathbf{p}_i^{DVL}\}_{i=1}^N$  and the camera-based posi-  
 282 tions  $\{\mathbf{p}_i^{CAM}\}_{i=1}^N$ , it is necessary to determine the optimal similarity transformation  $S^* = \{s^*, R_c^{b*}, \mathbf{t}_{c,b}^*\} =$   
 283  $\{s^*, \psi_c^{b*}, t_{c,b}^{x*}, t_{c,b}^{y*}\}$  that satisfies the minimization problem reported in Eq. 15.

$$284 \quad S^* = \underset{s, R_c^b, \mathbf{t}_{c,b}}{\operatorname{argmin}} \sum_{i=1}^N \|\mathbf{p}_i^{DVL} - sR_c^b \mathbf{p}_i^{CAM} - \mathbf{t}_{c,b}\|^2 \quad (15)$$

285 where it is necessary to suppose that

$$286 \quad R_c^b = R_z(\psi_c^b) \quad (16)$$

$$287 \quad \mathbf{t}_{c,b} = [t_{c,b}^x \ t_{c,b}^y \ 0]^\top \quad (17)$$

289 The solution of this least squares problem can be found using the method explained in (Umeyama, 1991).

290 The second step works with Ceres Solver, an open-source library that provides a rich set of tools to construct  
 291 and solve an optimization problem. Ceres solves robustified bounds constrained non-linear least squares  
 292 problems of the form:

$$293 \quad \min_{\mathbf{x}} \quad \frac{1}{2} \sum_i \rho_i (\|f_i(x_{i_1}, \dots, x_{i_k})\|^2). \quad (18)$$

$$l_j \leq x_j \leq u_j$$

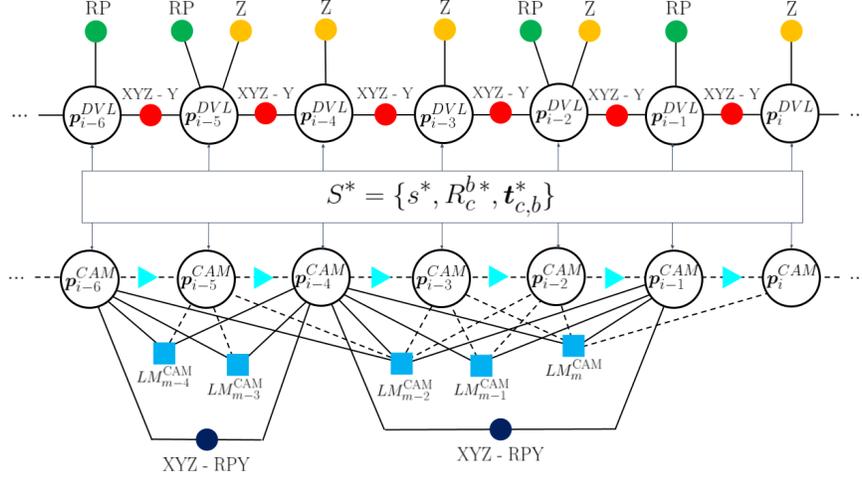


Figure 2: Comparison of the two factor graphs (e.g., the DVL-based graph on the top and the camera-based graph on the bottom of the image) employed for the scale factor ambiguity resolution. The dashed lines in the bottom graph are the edges which are not reported in the whole graph. For ease of reading, a one-to-one association between the two graphs is considered.

294 The expression  $\rho_i (\|f_i(x_{i_1}, \dots, x_{i_k})\|^2)$  represents the residual block, where  $\rho_i(\cdot)$  is the loss function used to  
 295 reduce the influence of outliers on the solution and  $f_i(\cdot)$  is the cost function that depends on the parameters  
 296 block  $\{x_{i_1}, \dots, x_{i_k}\}$ .  $l_j$  and  $u_j$  are the lower and upper bounds on the parameter block  $x_j$ .

297 Defining the state  $\mathbf{x} = [s \ \phi_c^b \ \theta_c^b \ \psi_c^b \ t_{c,b}^x \ t_{c,b}^y \ t_{c,b}^z]^\top$ , the loss function is assumed to be the identity  
 298 function, the cost function is the same as in the first step of the algorithm

$$299 \quad f(\mathbf{x}) = \mathbf{p}_i^{DVL} - s R_c^b \mathbf{p}_i^{CAM} - \mathbf{t}_{c,b} \quad (19)$$

300 where, unlike the previous case, it is supposed that

$$301 \quad R_c^b = R_z(\psi_c^b) R_y(\theta_c^b) R_z(\phi_c^b) \quad (20)$$

$$302 \quad 303 \quad \mathbf{t}_{c,b} = [t_{c,b}^x \ t_{c,b}^y \ t_{c,b}^z]^\top. \quad (21)$$

304 The initial guess and the upper and lower bounds are computed thanks to the values estimated in the closed-  
 305 form solution. Considering that this is a small problem with few parameters and relatively dense Jacobians,  
 306 dense QR factorization is the method of choice (Björck, 1996).

## 307 5.2 Reset procedures

308 Although iSAM2 reduces the variables to be optimized to a small subset, it is necessary to apply a reset  
 309 procedure to maintain a limited factor graph and avoid increasing nodes and edges. In particular, considering  
 310 that the presence of visual landmark nodes constrains several pose nodes, the computational burden tends  
 311 to increase at every iteration step, and the factor graph is more arduous to be managed. Two factor graph  
 312 reset procedures have been developed to avoid the increase of the graph size, where the first is dedicated to  
 313 compacting the factor graph without reducing the visual landmark nodes, and the second operates on the  
 314 whole factor graph reducing all the information to the ones contained in the last node. While the first reset  
 315 strategy will be called keyframe reset, the second one will be referred as global reset. One of the two reset  
 316 strategies is applied when the number of pose nodes of the factor graph reaches a value equal to  $N$ . The  
 317 status of the factor graph is checked to decide which one of the two strategies are applied. In particular, the  
 318 keyframe reset procedure is recalled only if the visual SLAM algorithm is active and for a maximum number  
 319 of consecutive times equal to  $p$ . The last condition is set to maintain control of the increase of the execution

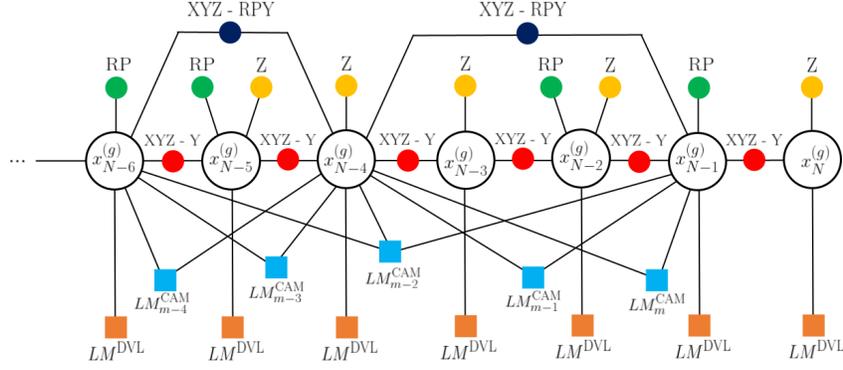


Figure 3: Last nodes of the factor graph  $g$  constrained with vision-based landmarks and all the onboard sensors.

320 time of each filter iteration. When the visual SLAM part of the navigation algorithm is not working due to  
 321 the external visibility conditions or when the factor graph is reset for the  $(p + 1)$ -th time, the global reset  
 322 algorithm is employed. It is necessary to notice that the keyframe reset procedure does not delete all the  
 323 information contained in the previous pose graph. Still, only the ones related to the IMU, DVL, and DS  
 324 measurements are removed. Indeed, this information is compressed in a new framework, which contains all  
 325 the properties to be transferred from the previous to the following factor graph. On the contrary, the global  
 326 factor reset reduces all the information to be transferred to the new factor graph to the ones in the last node  
 327 of the previous factor graph.

328 Both the reset strategies are now analyzed in detail to outline which information is passed from the previous  
 329 to the actual graph and how these measurements are compressed in the new framework. Considering the  
 330 keyframe reset procedure and referring to Fig. 3 and Fig. 4, the following actions are performed to obtain  
 331 the graph  $g + 1$  from the graph  $g$ .

- 332 • The  $i + 1$  keyframe pose nodes are transferred from the previous to the actual factor graph. The  
 333 first keyframe node, as the one associated with the state  $x_k^{(g)}$ , is constrained with a prior factor with  
 334 the last estimated value. All the subsequent  $i - 1$  keyframe nodes are determined by an XYZ-RPY  
 335 factor obtained from each last estimated value and the associated covariance.
- 336 • All the  $m + 1$  visual landmark points are transferred from the previous to the actual factor graph.  
 337 They are employed to maintain constraints between all the keyframe pose nodes. Each landmark  
 338 node is reported in the current graph with its last estimate and covariance and all the vision-based  
 339 edges.
- 340 • The last pose node associated with the state  $x_N^{(g)}$ , even if it is not a keyframe node, is transferred to  
 341 the actual graph to be employed as starting point to insert the acquired measurements as constraints.  
 342 This node is constrained to the last keyframe node with an XYZ-RPY odometry factor computed  
 343 from the last pose estimated values of the two nodes. The relative rototranslation transformation is  
 344 thus computed and applied as a constraint.

345 All the DVL-based landmarks are reported in the global NED reference frame using the poses estimated with  
 346 the graph  $g$ , and they are employed to build the point cloud for the seabed reconstruction. Even though the  
 347 whole graph has been reset, the visual SLAM part, if the visibility is acceptable, continues to compute poses  
 348 and visual landmarks, which are inserted in the new graph and connected to the keyframe nodes passed from  
 349 the previous graph. Furthermore, until a new keyframe is not computed, the new nodes are inserted thanks  
 350 to the DVL-based DR, the DS measurements, and the attitude estimator filter outputs.

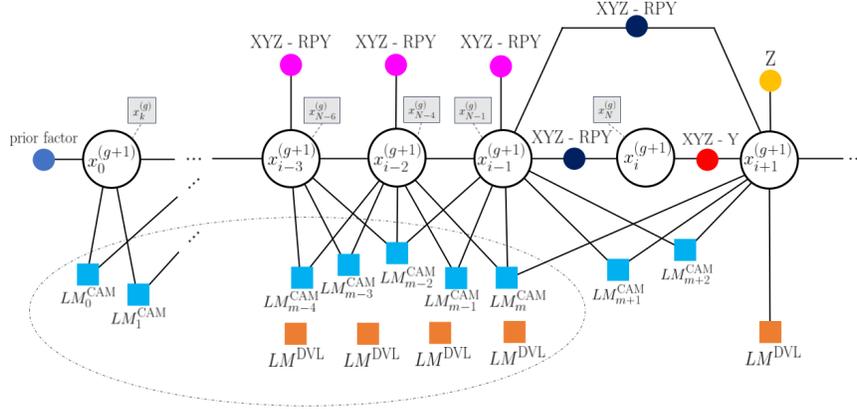


Figure 4: First nodes of the factor graph  $g + 1$  after the employment of the keyframe reset procedure. The values in the grey boxes represent the corresponding states taken from the previous factor graph  $g$  and transferred to the actual graph  $g + 1$ .

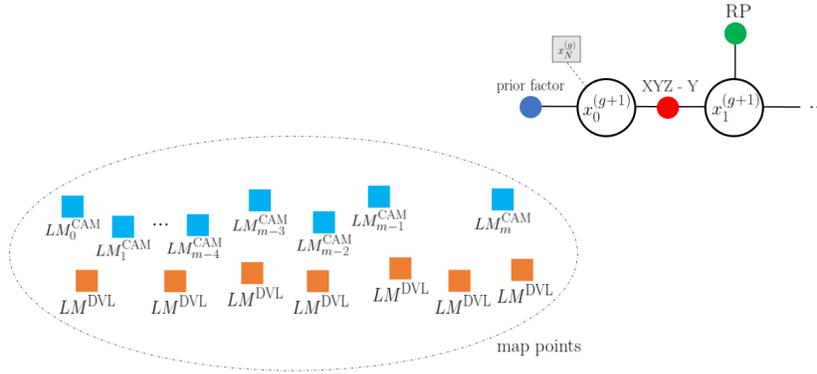


Figure 5: First nodes of the factor graph  $g + 1$  after the employment of the global reset procedure. The values in the grey boxes represent the corresponding states taken from the previous factor graph  $g$  and transferred to the actual graph  $g + 1$ .

351 Considering the global reset procedure and referring to Fig. 3 and Fig. 5, the following actions are performed  
 352 to obtain the graph  $g + 1$  from the graph  $g$ .

- 353 • Only the last pose node associated with the state  $x_N^{(g)}$  is transferred to the actual graph to be  
 354 employed as starting point to insert the acquired measurements as constraints. It is constrained  
 355 with a prior factor with the last estimated value.
- 356 • The visual landmarks and the keyframe poses are not transferred from the previous to the actual  
 357 graph. All positions of the estimated DVL-based and visual landmarks are saved as estimated in  
 358 the last optimization of the previous graph, and they are employed to build the point cloud for the  
 359 seabed reconstruction.

360 Even if the visibility is acceptable, the visual SLAM algorithm is reinitialized, the scale factor is again  
 361 computed, and no information is transferred from the vision-based part of the previous graph. Despite the  
 362 loss of some helpful information, the global reset procedure is necessary to limit the algorithm's computation  
 363 burden.

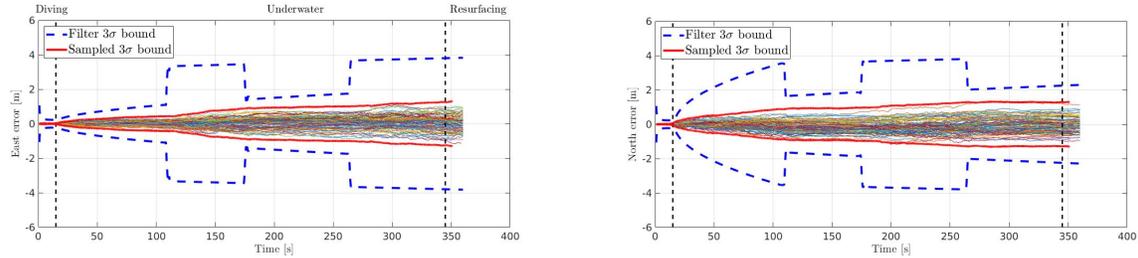


Figure 6: East and North position estimation errors versus their  $3\sigma$  bounds obtained from 100 simulation analysis with the SLAM algorithm. The  $\sigma$  values are computed as the square-root of the corresponding diagonal element of the estimated covariance matrix.

## 364 6 Navigation results in simulated environment

365 To validate the developed DVL and camera-based SLAM algorithm, realistic simulations were performed by  
 366 means of the Unmanned Underwater Vehicle Simulator (UUV Simulator). In particular, while navigation  
 367 performance has been evaluated thanks to a Monte Carlo simulation, mapping capabilities have been analyzed  
 368 with a lawnmower survey at a constant depth over a simulated seabed generated with a known mathematical  
 369 function  $z = f(x, y)$ . The obtained results have been employed to evaluate the goodness of the whole  
 370 algorithm and some of its main features, such as the reset procedure and the scale factor computation  
 371 algorithm. To focus attention on the navigation and mapping capabilities of the filter, the DVL and the  
 372 camera have been modeled thanks to the simulator features. The realistic simulations were based on the  
 373 dynamic model of FeelHippo AUV implemented in the UUV Simulator and on modeling all the onboard  
 374 sensors. In particular, the DVL beams have been modeled by applying a noise in the measured value, which  
 375 determines a noise in the measured velocity. The camera has been modeled with a noise in the pixel position  
 376 of the acquired image, which influences both the vehicle and landmark position estimation.

377 During the Monte Carlo simulations, the position filter was fed with the data coming from the simulated  
 378 sensors, as the GNSS, when the vehicle was higher than a fixed depth, depth sensor, DVL and camera. To  
 379 increase adherence to the real dataset, the DVL speed measurements have been published with a 5 Hz rate,  
 380 and the camera acquired images with a frequency of 10 Hz. The proposed strategies have been tested on a  
 381 vehicle whose dynamic behavior has been simulated using the model implemented in UUV Simulator, which  
 382 has traveled a rectangular path at a fixed depth of 2 m. A Monte Carlo simulation with 100 iterations has  
 383 been performed. The position errors and the estimated  $3\sigma$  bounds along the East and North directions are  
 384 reported in Fig. 6. The covariance trend follows the trajectory described by the vehicle. Still, the SLAM  
 385 algorithm, due to the presence of visual landmarks which constrains the vehicle position, provides an elliptic  
 386  $3\sigma$  bound with major axis perpendicular to the direction followed by the vehicle. Despite its particular shape,  
 387 the  $3\sigma$  bound continuously diverges when the vehicle is under the sea surface, and no position measurements  
 388 are available, correctly representing the behavior of the AUV.

389 Furthermore, as in the previous section, the estimated resurfacing position has been compared with the  
 390 theoretical first GNSS fix and its  $3\sigma$  bound. The resurfacing positions estimated in all the Monte Carlo  
 391 simulations fall inside the  $3\sigma$  bound, guaranteeing reasonable estimations. Furthermore, it is possible to  
 392 compare the  $3\sigma$  bound estimation obtained from the filter and the  $3\sigma$  bound estimation obtained from the  
 393 simulated data, evaluating the latter by computing the best normal distribution approximating the estimated  
 394 resurfacing positions with respect to the theoretical ones (Fig. 7).

395 Analyzing the results obtained from the lawnmower survey at a constant depth of 5 meters and comparing  
 396 the estimated trajectory with the ground truth provided by the simulator, it is possible to notice that  
 397 the divergence over time of the navigation error is reduced (see Fig. 8). Indeed, even if a global loop  
 398 closure on the visual keyframes is not performed, the presence of the highly accurate DVL measurements  
 399 can maintain a low estimation error drift. Furthermore, Fig. 8 shows the estimated trajectory on the NED  
 400 frame, where it is possible to notice the points where the system has been reset. Considering that the

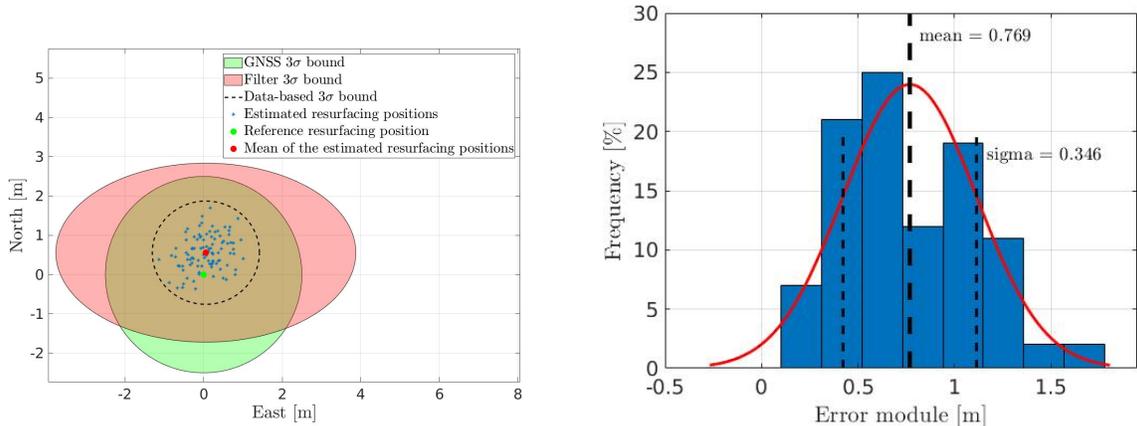


Figure 7: On the left, the estimated resurfacing positions versus the theoretical GNSS fix position obtained from 100 simulation analysis for the SLAM algorithm. On the right, histograms containing the estimated resurfacing position errors obtained from 100 simulation analysis for the SLAM algorithm.

Table 1: FeelHippo AUV Main Properties

Weight [kg]	35
Dimensions [mm]	600×640×500
Maximum Depth [m]	30
Maximum Longitudinal Speed [m/s]	1
Battery Life [h]	3

401 simulated seabed has been textured with a feature-rich image, it is necessary to see that the visual part  
 402 of the SLAM algorithm continues to work for the whole trajectory. Thus, both reset strategies have been  
 403 employed to limit the computational burden. Fig. 9 reports the estimated trajectory and the generated  
 404 point cloud. It is possible to evaluate the algorithm mapping capabilities by comparing the estimated point  
 405 cloud and the function employed to simulate the seabed. Considering that several outliers are kept in the  
 406 point cloud during the SLAM algorithm, which negatively influences the seabed reconstruction, the estimated  
 407 landmarks are elaborated to eliminate the wrong points and to downsample the cloud. Consequently, the  
 408 seabed reconstruction capabilities of the developed algorithm are analyzed in Section 8, where the employed  
 409 post-processing strategies are described.

## 410 7 Experimental results

411 The presented navigation and mapping strategy has been tested and validated by employing experimental  
 412 data recorded in Stromboli Island, Messina (Italy), in September 2022, during an autonomous underwater  
 413 mission performed in the framework of the project PATHFinder. During its autonomous navigation along  
 414 a pre-programmed path, the payload sensors were switched on, and the vehicle acquired both acoustic  
 415 and optical data. GNSS readings obtained from the satellites of the Galileo system were collected before  
 416 FeelHippo AUV (Allotta et al., 2017) dove and after it resurfaced. They have been employed as ground  
 417 truth to compute the resurfacing error and to globally reference the trajectory and the map.  
 418 FeelHippo AUV (see Fig. 10) is a compact vehicle capable of performing missions in shallow waters. The  
 419 main features of FeelHippo AUV are reported in Tab. 1. Furthermore, the sensors available on board are  
 420 listed as follows:

- 421 • U-blox 7P precision GNSS;

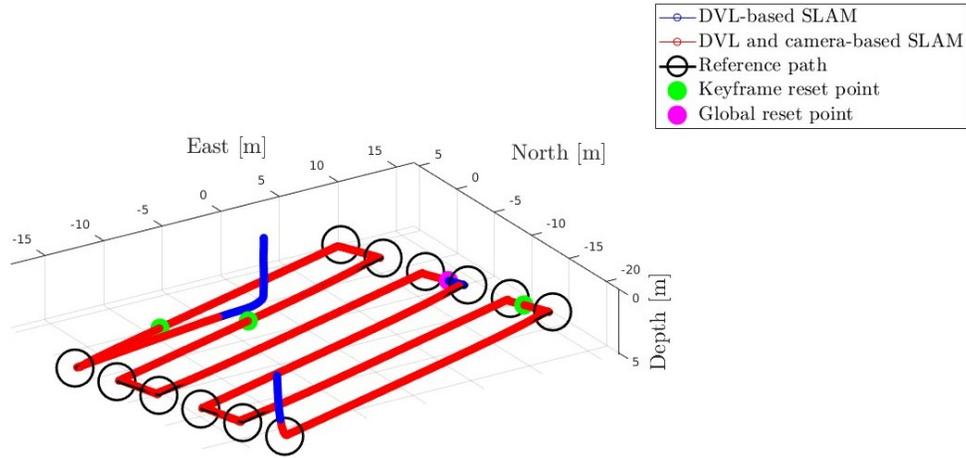


Figure 8: 3D plot of the estimated trajectory in the NED reference system, where the reset points and the areas where vision is not working are highlighted.

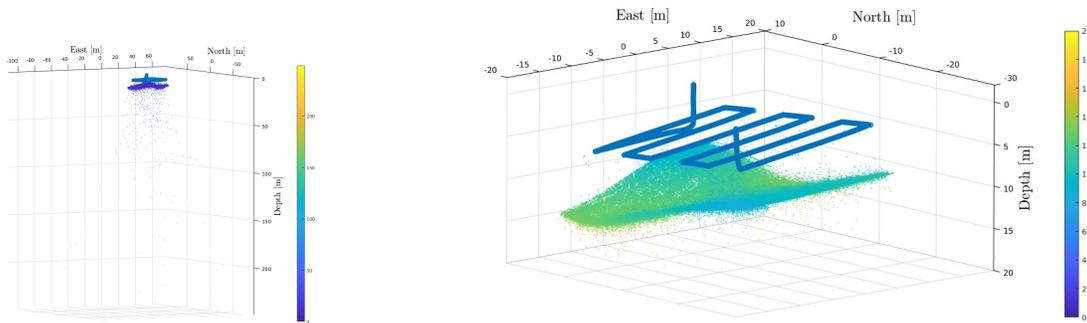


Figure 9: Representation of the point cloud and the travelled trajectory estimated through the SLAM algorithm. While on the top image the entire point cloud is reported and, due to the presence of outliers, the depth scale is too extended, on the bottom image a zoom on the region of interest is performed.



Figure 10: FeelHippo AUV before an on-field underwater mission.

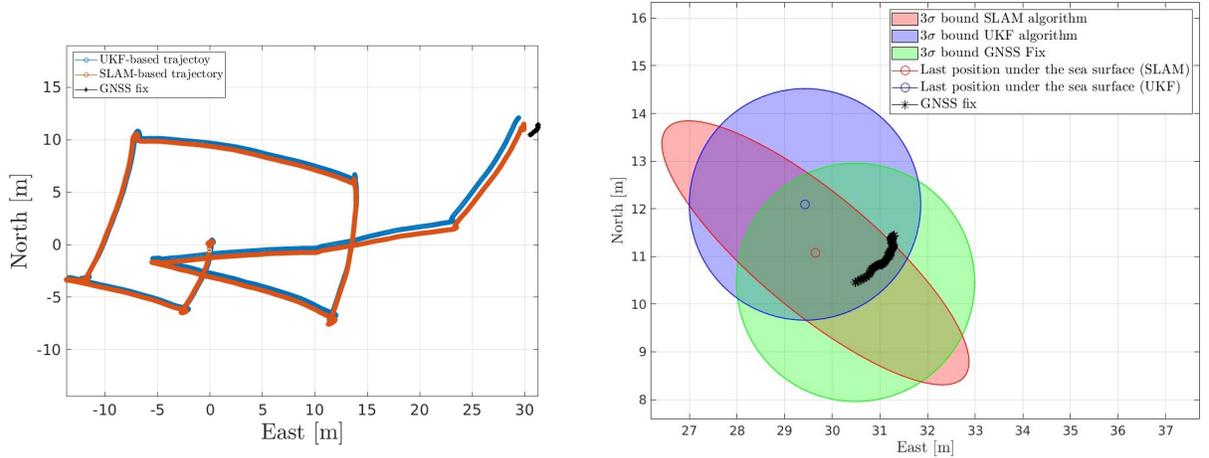


Figure 11: On the left, comparison between the trajectories estimated with the SLAM algorithm and the UKF algorithm during the mission accomplished in Stromboli Island, Messina (Italy). A ground truth, when the vehicle was under the sea surface, is not available, but the first GNSS fix when the vehicle resurfaces can be employed as reference to evaluate the resurfacing error. On the right,  $3\sigma$  bound of the last positions under the sea surface estimated with the SLAM and UKF algorithms and the first GNSS fix measurement with its accuracy  $3\sigma$  bound.

Table 2: Navigation performance for the mission accomplished in Stromboli Island, Messina (Italy): resurfacing error.

Navigation strategy	Error [m]
UKF algorithm	1.943
SLAM algorithm	0.899

- 422 • Orientus Advanced Navigation IMU;
- 423 • KVH DSP 1760 single-axis high precision Fiber Optic Gyroscope (FOG);
- 424 • Nortek DVL1000 DVL, measuring linear velocity and acting as DS;
- 425 • Teledyne BlueView M900 2D Forward Looking SONAR (FLS);
- 426 • two Microsoft Lifecam Cinema forward- and bottom-looking cameras.

427 The developed SLAM strategy has been compared with the Standard UKF algorithm chosen from the  
 428 navigation strategies proposed in (Bucci et al., 2023). The position resurfacing error values and covariances  
 429 have been evaluated on the North-East plane. Fig. 11 reports the estimated trajectories and an analysis of  
 430 the resurfacing errors with their  $3\sigma$  bound. From Tab. 2, analyzing the results from the GNSS resurfacing  
 431 error, it is easily noticeable that both the proposed strategies are acceptable in terms of navigation estimation  
 432 quality.

433 To evaluate the agreement between estimation errors and estimated uncertainty, the  $3\sigma$  bounds during the  
 434 resurfacing phase are presented. This is summarized in Fig. 11, where the  $3\sigma$  bounds for the filters and  
 435 the GNSS are presented. In all the analyzed cases, the position provided by the filter (with its confidence  
 436 bounds) appears to guarantee a reasonable prediction of the vehicle’s true position when it resurfaces. The  
 437 employed GNSS has an expected accuracy on the order of meters and the 2D error can be represented as a  
 438 2D Gaussian distribution whose components are independently distributed.

439 Focusing the attention on the SLAM algorithm and its mapping capabilities, Fig. 12 reports the SLAM-  
 440 based estimated trajectory and the generated point cloud. It is possible to evaluate the algorithm mapping

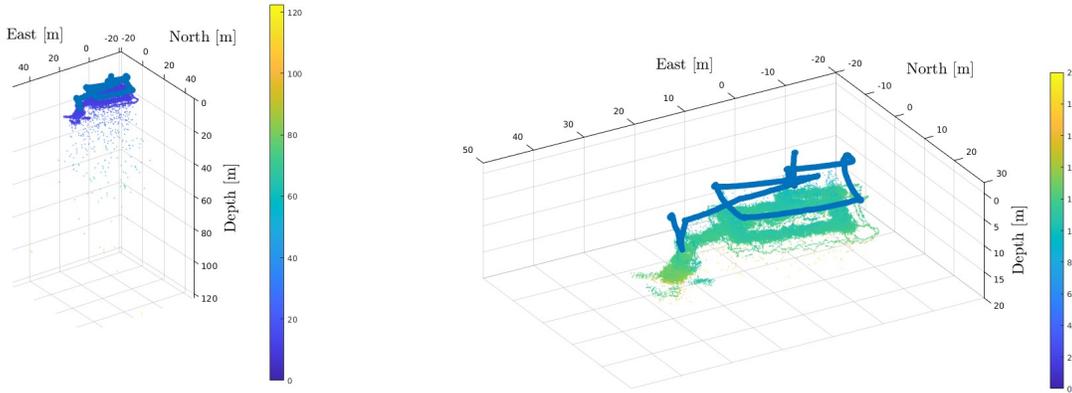


Figure 12: Representation of the point cloud and the travelled trajectory estimated through the SLAM algorithm during the mission in Stromboli Island, Messina (Italy). While on the top image the entire point cloud is reported and, due to the presence of outliers, the depth scale is too extended, on the bottom image a zoom on the region of interest is performed.

Table 3: Estimated scale factor and rototranslation transform between DVL and camera reference systems.

Parameter	Initial guess	Value after step 1	Value after step 2
$s$	0.0	5.448	5.529
$\phi_c^b$ [deg]	0.0	0.0	-0.005
$\theta_c^b$ [deg]	90.0	90.0	89.477
$\psi_c^b$ [deg]	0.0	10.119	8.43
$t_{c,b}^x$ [m]	0.24	0.24	0.233
$t_{c,b}^y$ [m]	0.07	0.06	0.076
$t_{c,b}^z$ [m]	0.05	0.05	0.049

441 capabilities by comparing the estimated point cloud with a bathymetry of the region around the island. As  
 442 for the test in simulated environments, several outliers are kept in the point cloud during the SLAM algorithm  
 443 operation, which negatively influences the seabed reconstruction. Consequently, the seabed reconstruction  
 444 capabilities of the developed algorithm and the comparison with the ground truth bathymetry are reported  
 445 in Section 8, where the employed post-processing strategies are described.

446 The scale factor computation procedure has been applied to estimate the scale factor between the DVL-  
 447 based trajectory and the visual part of the algorithm before fusing them in the whole factor graph. In  
 448 particular, approximate values of the relative position and orientation between the DVL and the camera has  
 449 been provided as input to the algorithm, but their values have been kept as variables in the optimization  
 450 process. The scale factor between the DVL-based trajectory and the visual SLAM has been solved with the  
 451 developed algorithm, and the results have been reported in Tab. 3. It is necessary to highlight that the  
 452 proposed strategy can compensate the alignment error between the camera and the DVL frames. Indeed, due  
 453 to uncontrollable external conditions (e.g., loosening of the screws during the vehicle preparation, collisions  
 454 during the diving procedure), the camera rotated around its  $z$ -axis during the autonomous mission of an  
 455 unknown quantity which has been estimated and compensated by the algorithm. The resurfacing error value  
 456 is equal to 0.899 meters, indicating a high navigation accuracy of the proposed strategy with respect to the  
 457 GNSS fixes obtained when the vehicle resurfaced.

458 Finally, regarding the computational burden, the execution time of the filter has been subject of the analysis.  
 459 The sum of the requested time to perform the measurement insertion in the factor graph and the optimization

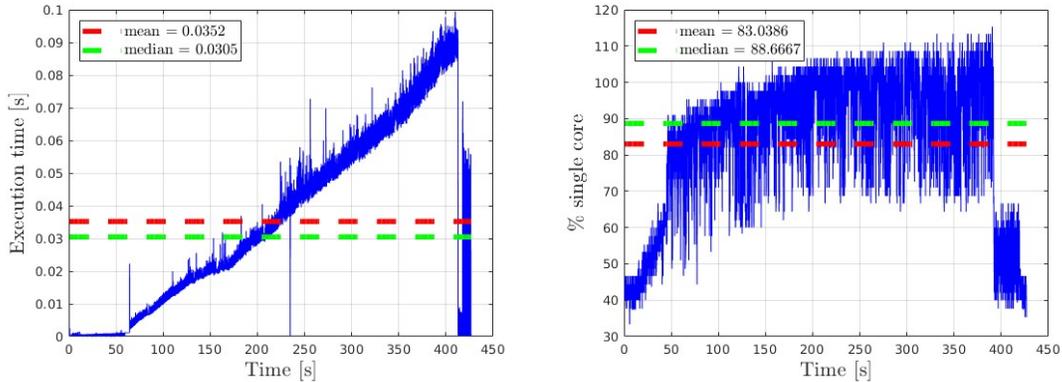


Figure 13: On the left, execution time of the SLAM filter, calculated at each iteration as the sum of the requested time for measurement insertion in the factor graph and for the optimization process. On the right, Central Processing Unit (CPU) burden analysis. In red and green are respectively reported the mean and the median.

Table 4: Mean with the associated covariance and median values of the depth error in presence and in absence of the filtering procedure.

Parameter	Before filtering	After filtering
Mean [m]	0.2767	0.2002
Covariance [m]	7.6212	0.0386
Median [m]	0.1469	0.1465

460 process has been considered. For what concern the CPU analysis, the output of the command *top* has been  
 461 recorded to store the data. The results can be found in Fig. 13. It is necessary to notice that the instants  
 462 where the visual part of the algorithm is initialized and stopped can be easily highlighted thanks to its  
 463 influence on the execution time of each iteration. Indeed, despite the SLAM algorithm optimizes only the  
 464 last nodes thanks to the properties of the iSAM2 library, handling a continuously growing point cloud  
 465 increases the required computational cost. When the vehicle resurfaces and the visual part of the algorithm  
 466 is excluded due to visibility limitations, the necessary computational burden drastically decreases. Indeed,  
 467 the point cloud is saved, and only the position nodes are updated when new measurements are acquired.

## 468 8 Mapping performance analysis

469 Mapping the surrounding environment is a common task in underwater exploration, and it is fundamental  
 470 to enhance the vehicle capabilities to find objects of potential interest. The point clouds obtained from  
 471 the SLAM algorithm have been processed with an automatic tool to obtain a 3D reconstruction of the sea  
 472 bottom. The developed reconstruction strategy takes as input the estimated point cloud and the geographi-  
 473 cal coordinate of a reference point and automatically generates a 3D reconstruction and a georeferenced  
 474 depth map, thanks to the employment of the functions implemented in the open-source libraries Point Cloud  
 475 Library (PCL) (Rusu and Cousins, 2011) and Open3d (Zhou et al., 2018).

476 Analyzing the point cloud obtained from the navigation algorithm applied in both simulated and real envi-  
 477 ronments, it is necessary to notice that some points can be classified as outliers. Therefore the need arises to  
 478 eliminate them as the displayed graphs are excessively bulky and negatively influence the mesh realization.  
 479 For each point, a fixed number of neighbors is defined to estimate the mean of the average distance, and a  
 480 point is considered an outlier if the average distance to its neighbors is above a specified threshold (Rusu  
 481 et al., 2008). The outlier eliminating process, therefore, leads to a significant decrease in points, making the  
 482 representations more uniform. Subsequently, the point cloud is processed with a smoothing method to filter

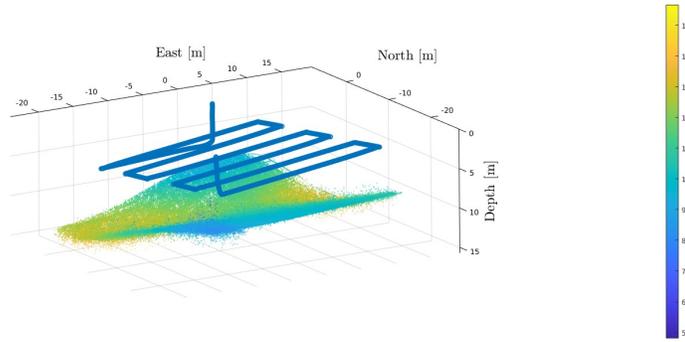


Figure 14: Filtered point cloud and estimated trajectory travelled by the simulated vehicle.

483 out the noise of the measurements on the processed points. In particular, median filtering of the 3D point  
 484 cloud data is performed.

485 The 3D mesh generation algorithm is applied to the filtered point cloud. The Poisson algorithm (Kazhdan  
 486 et al., 2006) is applied, and its parameters have been set to optimize the reconstruction process. It is  
 487 necessary to note that the depth value and the limit density of the points at which the reconstruction is  
 488 cut have been chosen to compromise reconstruction speed and estimation quality. The advantages of the  
 489 Poisson algorithm over other surface fitting methods are numerous. Many implicit methods of surface fitting  
 490 segment the data into regions for local fitting and further combine these local approximations using blending  
 491 functions. In contrast, Poisson reconstruction is a global solution that considers all the data simultaneously  
 492 without resorting to heuristic partitioning or blending. In this way, Poisson reconstruction creates very  
 493 smooth surfaces that robustly approximate noisy data.

494 Firstly, considering that in the simulated environment created with UUV Simulator the seabed can be  
 495 generated with a mathematical function  $z = f(x, y)$ , it is possible to evaluate the performance of the filtering  
 496 algorithm. It is also necessary to notice that the simulated seabed has been textured with an image rich  
 497 in features to facilitate the correct behavior of the visual part of the SLAM algorithm. Fig. 14 reports the  
 498 3D filtered point cloud with the estimated trajectory. It is necessary to compare this point cloud with the  
 499 one directly obtained from the SLAM algorithm and reported in Fig. 9. Two error maps have been created  
 500 with the point clouds, as before and after the filtering procedure, to analyze the improvements in seabed  
 501 reconstruction. It is necessary to notice that the outlier points are correctly removed, and the point cloud  
 502 size is reduced to increase its easiness of management by the reconstruction algorithm (see Fig. 15 and Fig.  
 503 16). As can be retrieved from Tab. 4, while the outlier removal process does not influence the mean and the  
 504 median values due to the high number of points, the covariance associated with the mean value is strongly  
 505 reduced.

506 Finally, the 3D point cloud has been processed with the Poisson reconstruction algorithm to build the  
 507 3D mesh. Thanks to the chosen reconstruction algorithm, the obtained mesh is smoothed and correctly  
 508 follows the shape of the simulated sea bottom (see Fig. 17). Some of the results obtained during the  
 509 mission performed in Stromboli Island, Messina (Italy), in September 2022, are presented. In particular, a  
 510 3D reconstruction and a geo-localized map of the sea bottom are reported. The reconstruction comprises  
 511 around 240k points obtained as output from the SLAM algorithm. Firstly, the outlier points have been  
 512 removed (see Fig. 19), and then, the 3D point cloud has been processed with the Poisson reconstruction  
 513 algorithm to build the 3D mesh, which is shown in Fig. 20.

514 The good matching between the reference bathymetry, whose data have mainly funded by the National  
 515 Research Council and Presidenza del Consiglio dei Ministri–Dipartimento della Protezione Civile, through  
 516 specific agreement (see Fig. 18), and the estimated 3D reconstruction can also be observed to prove the

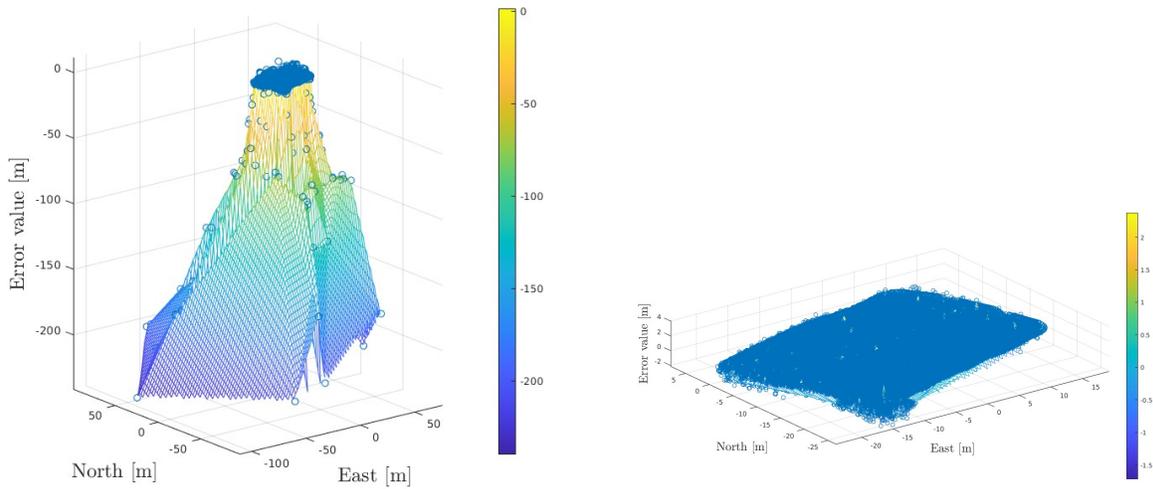


Figure 15: Representation of the error point clouds computed by comparing the reference sea bed function and the estimated point cloud and generation of the estimated error maps before (top image) and after (bottom image) the filtering procedure.

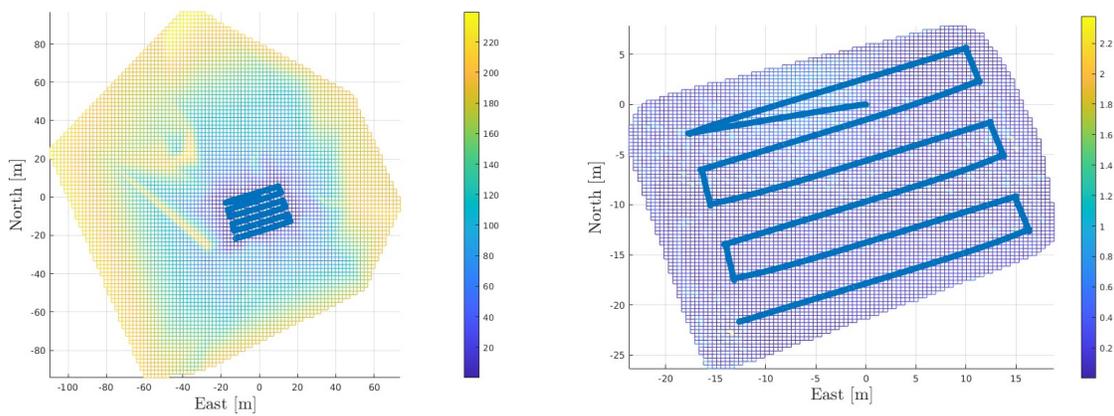


Figure 16: Comparison between the estimated error maps before (top image) and after (bottom image) the filtering procedure with respect to the travelled trajectory by the simulated vehicle.

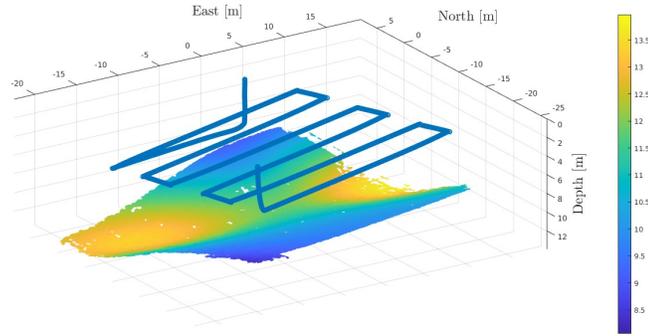


Figure 17: Resulting sea bottom 3D mesh reconstruction and estimated trajectory travelled by the simulated vehicle.

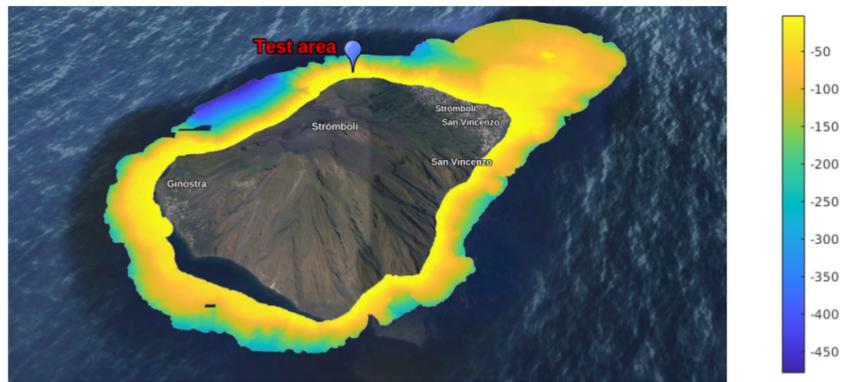


Figure 18: Reference bathymetry of the sea bottom around Stromboli Island, Messina (Italy). The test area, where FeelHippo AUV performed its autonomous mission, is highlighted.

517 reconstruction’s goodness. The provided bathymetry has a horizontal resolution of 5 meters. Thus only  
 518 an approximate comparison can be performed, but it can be sufficient to have a simple evaluation of the  
 519 generated point cloud. All the points of the cloud that lies in each square generated from the ground truth  
 520 bathymetry are employed to compute the mean point and perform the comparison (Fig. 21).

## 521 9 Conclusion and future developments

522 Considering that Kalman filtering condenses the vehicle’s history in the last estimate and covariance, a  
 523 MAP strategy based on factor graphs has been developed to overcome these limitations and include visual  
 524 landmarks in the estimation process. Visual features are sometimes difficult to be found in underwater  
 525 environments due to visibility and texture issues. Consequently, the strategy fuses DVL measurements with  
 526 a visual SLAM system to simultaneously perform accurate navigation and mapping tasks. DVL beam data  
 527 can be employed for speed measurement and to obtain an approximated knowledge of the sea bottom. Both  
 528 simulated and experimental data have been employed to evaluate the capabilities of the developed strategy.  
 529 The experimental data have been acquired during trials at Stromboli Island, Messina (Italy).  
 530 During the experimental campaign, FeelHippo AUV was the only vehicle involved; nevertheless, since the  
 531 proposed solution is not tailored to a particular vehicle, its outcomes can be deemed as general, and future  
 532 developments will include the testing of the navigation strategy on other vehicles. Furthermore, progresses  
 533 on the developed algorithms still needs to be made. Integrating the developed estimation strategy within

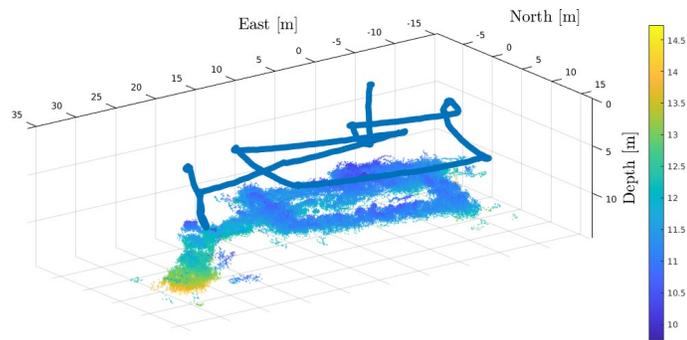


Figure 19: Filtered point cloud and estimated trajectory travelled by the vehicle during the autonomous mission accomplished in Stromboli Island, Messina (Italy).

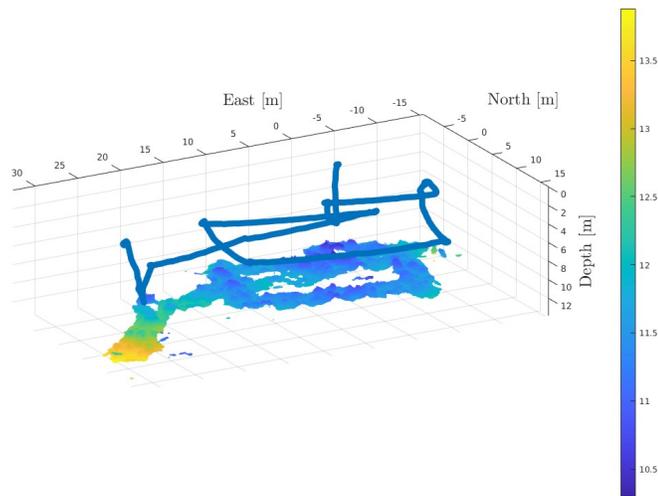


Figure 20: Resulting sea bottom 3D mesh reconstruction and estimated trajectory travelled by the vehicle during the autonomous mission accomplished in Stromboli Island, Messina (Italy).

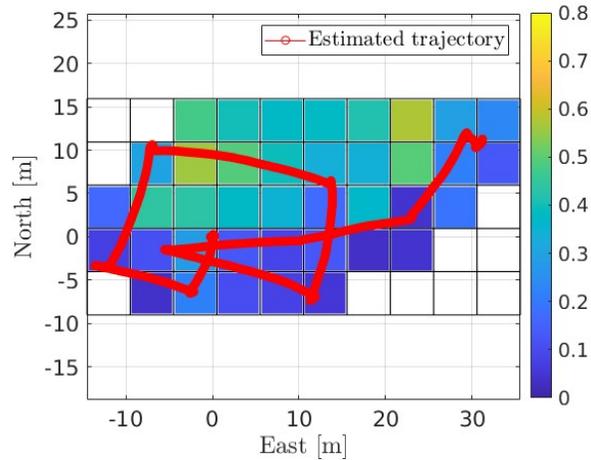


Figure 21: Estimated error bathymetry map with respect to the trajectory travelled by the vehicle during the autonomous mission accomplished in Stromboli Island, Messina (Italy).

534 the attitude estimator could represent an important subject to be investigated. Indeed, developing a unique  
 535 filter that works on both the attitude and position estimation in a coupled way could increase navigation  
 536 precision. Concerning the strategies for graph edge computation, including constraints obtained from acoustic  
 537 FLS images, which have been employed for speed estimation in (Bucci et al., 2023), in the pose graph  
 538 framework would push forward the performance of the navigation filter. Despite the intrinsic characteristic  
 539 (low resolution, influence of the viewpoint) of FLS images poses relevant issues to face, the employment  
 540 of an additional acoustic sensor can be useful to apply the developed strategy in scenarios with reduced  
 541 visibility. Finally, from a mapping-based point of view, a multi-vehicle solution for autonomously fusing  
 542 the underwater environment reconstructions could represent a coherent continuation of the research activity  
 543 carried out so far. The proposed SLAM strategy could operate onboard of each vehicle and, by employing  
 544 relative or absolute position measurements, the estimated maps could be fused in a unique more detailed  
 545 reconstruction.

## 546 10 Acknowledgement

547 The activities leading to the presented results have been carried out in the context of the PATHfinder  
 548 European project, part of the ESA’s NAVISP Programme, and with the support of the National Research  
 549 Council and Presidenza del Consiglio dei Ministri–Dipartimento della Protezione Civile.

## 550 References

- 551 Allotta, B., Conti, R., Costanzi, R., Fanelli, F., Gelli, J., Meli, E., Monni, N., Ridolfi, A., and Rindi,  
 552 A. (2017). A low cost autonomous underwater vehicle for patrolling and monitoring. *Institution of*  
 553 *Mechanical Engineers, Part M: Journal of Engineering for the Maritime Environment*, 231(3):740–749.
- 554 Björck, A. (1996). *Numerical Methods for Least Squares Problems*. Society for Industrial and Applied  
 555 Mathematics.
- 556 Bucci, A., Franchi, M., Ridolfi, A., Secciani, N., and Allotta, B. (2023). Evaluation of UKF-Based Fusion  
 557 Strategies for Autonomous Underwater Vehicles Multisensor Navigation. *IEEE Journal of Oceanic*  
 558 *Engineering*, 48(1):1–26.
- 559 Bucci, A., Ridolfi, A., Franchi, M., and Allotta, B. (2021). Covariance and gain-based federated unscented

560 kalman filter for acoustic-visual-inertial underwater navigation. In *OCEANS 2021: San Diego – Porto*,  
561 pages 1–7.

562 Bucci, A., Zacchini, L., Franchi, M., Ridolfi, A., and Allotta, B. (2022). Comparison of feature detection  
563 and outlier removal strategies in a mono visual odometry algorithm for underwater navigation. *Applied*  
564 *Ocean Research*, 118:102961.

565 Cadena, C., Carlone, L., Carrillo, H., Latif, Y., Scaramuzza, D., Neira, J., Reid, I., and Leonard, J. J. (2016).  
566 Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age.  
567 *IEEE Transactions on robotics*, 32(6):1309–1332.

568 Campos, C., Elvira, R., Rodríguez, J. J. G., M. Montiel, J. M., and D. Tardós, J. (2021). ORB-SLAM3: An  
569 Accurate Open-Source Library for Visual, Visual–Inertial, and Multimap SLAM. *IEEE Transactions*  
570 *on Robotics*, 37(6):1874–1890.

571 Cashmore, M., Fox, M., Larkworthy, T., Long, D., and Magazzeni, D. (2014). AUV mission control via  
572 temporal planning. In *2014 IEEE International Conference on Robotics and Automation (ICRA)*, pages  
573 6535–6541.

574 Castellanos, J. A., Neira, J., and Tardós, J. D. (2004). Limits to the consistency of ekf-based slam. *IFAC*  
575 *Proceedings Volumes*, 37(8):716–721. IFAC/EURON Symposium on Intelligent Autonomous Vehicles,  
576 Lisbon, Portugal, 5-7 July 2004.

577 Dellaert, F. (2012). Factor graphs and GTSAM: A hands-on introduction. Technical report, Georgia Institute  
578 of Technology.

579 Dellaert, F. (2021). Factor graphs: Exploiting structure in robotics. *Annual Review of Control, Robotics,*  
580 *and Autonomous Systems*, 4(1):141–166.

581 Dellaert, F. and Kaess, M. (2006). Square root sam: Simultaneous localization and mapping via square root  
582 information smoothing. *The International Journal of Robotics Research*, 25(12):1181–1203.

583 Dellaert, F. and Kaess, M. (2017). Factor graphs for robot perception. *Foundations and Trends in Robotics*,  
584 6(1-2):1–139.

585 Dissanayake, M., Newman, P., Clark, S., Durrant-Whyte, H., and Csorba, M. (2001). A solution to the simul-  
586 taneous localization and map building (slam) problem. *IEEE Transactions on Robotics and Automation*,  
587 17(3):229–241.

588 Du, P., Han, J., Wang, J., Wang, G., Jing, D., Wang, X., and Qu, F. (2017). View-based underwater SLAM  
589 using a stereo camera. In *OCEANS 2017 - Aberdeen*, pages 1–6.

590 Fallon, M. F., Folkesson, J., McClelland, H., and Leonard, J. J. (2013). Relocating Underwater Features  
591 Autonomously Using Sonar-Based SLAM. *IEEE Journal of Oceanic Engineering*, 38(3):500–513.

592 Ferri, G., Ferreira, F., and Djapic, V. (2017). Multi-domain robotics competitions: The CMRE experience  
593 from SAUC-E to the European Robotics League Emergency Robots. In *OCEANS 2017 - Aberdeen*,  
594 pages 1–7.

595 Forster, C., Carlone, L., Dellaert, F., and Scaramuzza, D. (2016). On-Manifold Preintegration for Real-Time  
596 Visual–Inertial Odometry. *IEEE Transactions on Robotics*, 33(1):1–21.

597 Franchi, M., Bucci, A., Zacchini, L., Ridolfi, A., Bresciani, M., Peralta, G., and Costanzi, R. (2021). Max-  
598 imum A Posteriori estimation for AUV localization with USBL measurements. *IFAC-PapersOnLine*,  
599 54(16):307–313. 13th IFAC Conference on Control Applications in Marine Systems, Robotics, and  
600 Vehicles CAMS 2021.

601 Grisetti, G., Guadagnino, T., Aloise, I., Colosi, M., Della Corte, B., and Schlegel, D. (2020). Least squares  
602 optimization: From theory to practice. *Robotics*, 9(3):51.

- 603 Hong, S. and Kim, J. (2020). Three-dimensional Visual Mapping of Underwater Ship Hull Surface Using  
604 Piecewise-planar SLAM. *International Journal of Control, Automation and Systems*, 18:564–574.
- 605 Huang, S. and Dissanayake, G. (2007). Convergence and consistency analysis for extended kalman filter  
606 based slam. *IEEE Transactions on Robotics*, 23(5):1036–1049.
- 607 Huang, T. A. and Kaess, M. (2015). Towards acoustic structure from motion for imaging sonar. In *2015*  
608 *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 758–765. IEEE.
- 609 Julier, S. and Uhlmann, J. (2001). A counter example to the theory of simultaneous localization and map  
610 building. In *Proceedings 2001 ICRA. IEEE International Conference on Robotics and Automation (Cat.*  
611 *No.01CH37164)*, volume 4, pages 4238–4243.
- 612 Kaess, M., Johannsson, H., Roberts, R., Ila, V., Leonard, J. J., and Dellaert, F. (2012). iSAM2: Incre-  
613 mental smoothing and mapping using the Bayes tree. *The International Journal of Robotics Research*,  
614 31(2):216–235.
- 615 Kaess, M., Ranganathan, A., and Dellaert, F. (2008). iSAM: Incremental smoothing and mapping. *IEEE*  
616 *Transactions on Robotics*, 24(6):1365–1378.
- 617 Kazhdan, M., Bolitho, M., and Hoppe, H. (2006). Poisson Surface Reconstruction. In *Proceedings of the*  
618 *Fourth Eurographics Symposium on Geometry Processing*, page 61–70.
- 619 Kim, A. and Eustice, R. M. (2013). Real-Time Visual SLAM for Autonomous Underwater Hull Inspection  
620 Using Visual Saliency. *IEEE Transactions on Robotics*, 29(3):719–733.
- 621 Mur-Artal, R., Montiel, J. M. M., and Tardós, J. D. (2015). ORB-SLAM: A Versatile and Accurate Monoc-  
622 ular SLAM System. *IEEE Transactions on Robotics*, 31(5):1147–1163.
- 623 Mur-Artal, R. and Tardós, J. D. (2017). ORB-SLAM2: An Open-Source SLAM System for Monocular,  
624 Stereo, and RGB-D Cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262.
- 625 Ozog, P. and Eustice, R. M. (2013). Real-time SLAM with piecewise-planar surface models and sparse 3D  
626 point clouds. In *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages  
627 1042–1049.
- 628 Paull, L., Saeedi, S., Seto, M., and Li, H. (2012). Sensor driven online coverage planning for autonomous  
629 underwater vehicles. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*,  
630 pages 2875–2880.
- 631 Prats, M., Ribas, D., Palomeras, N., García, J. C., Nannen, V., Wirth, S., Fernández, J. J., Beltrán, J. P.,  
632 Campos, R., Ridao, P., Sanz, P. J., Oliver, G., Carreras, M., Gracias, N., Marín, R., and Ortiz, A.  
633 (2012). Reconfigurable AUV for intervention missions: a case study on underwater object recovery.  
634 *Intelligent Service Robotics*, 5:19–31.
- 635 Rahman, S., Li, A. Q., and Rekleitis, I. (2018). Sonar Visual Inertial SLAM of Underwater Structures. In  
636 *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 5190–5196.
- 637 Rahman, S., Li, A. Q., and Rekleitis, I. M. (2018). SVIn2: Sonar Visual-Inertial SLAM with Loop Closure  
638 for Underwater Navigation. *ArXiv*, abs/1810.03200.
- 639 Rusu, R. B. and Cousins, S. (2011). 3D is here: Point Cloud Library (PCL). In *IEEE International*  
640 *Conference on Robotics and Automation (ICRA)*.
- 641 Rusu, R. B., Marton, Z. C., Blodow, N., Dolha, M., and Beetz, M. (2008). Towards 3D Point cloud based  
642 object maps for household environments. *Robotics and Autonomous Systems*, 56(11):927–941.
- 643 Umeyama, S. (1991). Least-squares estimation of transformation parameters between two point patterns.  
644 *IEEE Computer Architecture Letters*, 13(04):376–380.

- 645 Vidal, E., Palomeras, N., Istenič, K., Gracias, N., and Carreras, M. (2020). Multisensor online 3D view  
646 planning for autonomous underwater exploration. *Journal of Field Robotics*, 37(6):1123–1147.
- 647 Westman, E. and Kaess, M. (2018). Underwater AprilTag SLAM and calibration for high precision robot  
648 localization. Technical Report CMU-RI-TR-18-43, Carnegie Mellon University, Pittsburgh, PA.
- 649 Westman, E. and Kaess, M. (2019). Degeneracy-aware imaging sonar simultaneous localization and mapping.  
650 *IEEE Journal of Oceanic Engineering*, 45(4):1280–1294.
- 651 Westman, E. and Kaess, M. (2020). Degeneracy-Aware Imaging Sonar Simultaneous Localization and Map-  
652 ping. *IEEE Journal of Oceanic Engineering*, 45(4):1280–1294.
- 653 Youakim, D., Cieslak, P., Dornbush, A., Palomer, A., Ridao, P., and Likhachev, M. (2020). Multirepresenta-  
654 tion, Multiheuristic A\* search-based motion planning for a free-floating underwater vehicle-manipulator  
655 system in unknown environment. *Journal of Field Robotics*, 37(6):925–950.
- 656 Zacchini, L., Bucci, A., Franchi, M., Costanzi, R., and Ridolfi, A. (2019). Mono visual odometry for  
657 Autonomous Underwater Vehicles navigation. In *2019 MTS/IEEE Oceans, Marseille, France*.
- 658 Zhang, Y., Zhang, T., and Huang, S. (2018). Comparison of ekf based slam and optimization based slam  
659 algorithms. In *2018 13th IEEE Conference on Industrial Electronics and Applications (ICIEA)*, pages  
660 1308–1313.
- 661 Zhou, Q. Y., Park, J., and Koltun, V. (2018). Open3D: A modern library for 3D data processing.  
662 *arXiv:1801.09847*.