



**HAL**  
open science

## **AQUALOC: An Underwater Dataset for Visual-Inertial-Pressure Localization.**

Maxime Ferrera, Vincent Creuze, Julien Moras, Pauline Trouvé-Peloux

► **To cite this version:**

Maxime Ferrera, Vincent Creuze, Julien Moras, Pauline Trouvé-Peloux. AQUALOC: An Underwater Dataset for Visual-Inertial-Pressure Localization.. The International Journal of Robotics Research, 2019, 38 (14), pp.1549-1559. 10.1177/0278364919883346 . hal-02332498

**HAL Id: hal-02332498**

**<https://hal.science/hal-02332498>**

Submitted on 31 Oct 2019

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# AQUALOC: An Underwater Dataset for Visual-Inertial-Pressure Localization.

Journal Title  
XX(X):1-9  
©The Author(s) 2019  
Reprints and permission:  
sagepub.co.uk/journalsPermissions.nav  
DOI: 10.1177/ToBeAssigned  
www.sagepub.com/

SAGE

Maxime Ferrera<sup>1,2</sup>, Vincent Creuze<sup>2</sup>, Julien Moras<sup>1</sup> and Pauline Trouvé-Peloux<sup>1</sup>

## Abstract

We present a new dataset, dedicated to the development of simultaneous localization and mapping methods for underwater vehicles navigating close to the seabed. The data sequences composing this dataset are recorded in three different environments: a harbor at a depth of a few meters, a first archaeological site at a depth of 270 meters and a second site at a depth of 380 meters. The data acquisition is performed using Remotely Operated Vehicles equipped with a monocular monochromatic camera, a low-cost inertial measurement unit, a pressure sensor and a computing unit, all embedded in a single enclosure. The sensors' measurements are recorded synchronously on the computing unit and seventeen sequences have been created from all the acquired data. These sequences are made available in the form of ROS bags and as raw data. For each sequence, a trajectory has also been computed offline using a Structure-from-Motion library in order to allow the comparison with real-time localization methods. With the release of this dataset, we wish to provide data difficult to acquire and to encourage the development of vision-based localization methods dedicated to the underwater environment. The dataset can be downloaded from: <http://www.lirmm.fr/aqualoc/>

## Keywords

Dataset, Underwater robotics, Monocular Vision, IMU, Pressure, SLAM

## 1 Introduction

Accurate localization is critical for mobile robotics. In open outdoor areas, it can be obtained from Global Positioning System (GPS). However, in GPS-denied environments, such as indoor or beneath the sea surface, robots' position must be estimated from other sensors.

In underwater robotics, the localization problem is often solved by coupling high-grade Inertial Measurement Units (IMU) with compass, Doppler Velocity Logs (DVL) and pressure sensors (Paull et al. (2014)). Such solutions, classified as dead-reckoning (DR) localization, are highly dependent of the sensors quality and suffer from unbounded drift. While these methods can be employed quite safely for vehicles navigating in the middle of the water column (*i.e.* in obstacle free areas), they are not accurate enough for navigation in cluttered areas. In such places, Simultaneous Localization And Mapping (SLAM) methods are preferred. SLAM requires exteroceptive sensors, such as Lidar, sonar or camera, to measure the 3D structure of the environment. From these data, the localization is estimated while a 3D map is progressively built.

Visual SLAM (VSLAM) and Visual-Inertial Odometry (VIO) have been a hot research topic during the past decades (Cadena et al. (2016)). VSLAM consists in estimating localization from visual data, possibly enhanced by complementary sensors, through the mapping of the observed scenes. In ground and aerial robotics, the availability of many public datasets, such as KITTI (Geiger et al. (2012)), Malaga (Blanco et al. (2014)) or EuRoC (Burri et al. (2016)), to cite a few, has greatly impacted the development of VSLAM methods. Recent

algorithms, relying on monocular cameras (Mur-Artal et al. (2015); Forster et al. (2017); Engel et al. (2018)) or on visual-inertial sensors (Leutenegger et al. (2015); Mur-Artal and Tardos (2017); Qin et al. (2018)), have shown impressive results, with centimetric localization accuracy. In underwater robotics, many operations occur near the seabed (biology, Oil&Gas Industry, mine warfare, archaeology...), making visual information available. Nonetheless, in such conditions, the acquired images suffer from degradation like turbidity, backscattering and illumination issues, due to the medium properties. These poor imaging conditions must be accounted for in the development of underwater VSLAM or VIO systems, thus preventing use of the previously cited algorithms (Quattrini Li et al. (2017); Weidner et al. (2017); Zhang et al. (2018)). Some previous works have investigated the use of monocular camera for underwater localization (Burguera et al. (2015); Ferrera et al. (2019)), sometimes coupled to low-cost IMU and pressure sensor (Shkurti et al. (2011); Creuze (2017)), sonars (Rahman et al. (2018)) or even as a mean of detecting loop-closures in DR systems (Kim and Eustice (2013)). However, the limited amount of public datasets dedicated to this localization challenge prevent a fair comparison of these methods on common data. Moreover, the fact that these data are difficult to acquire, because of the required equipment and logistic,

<sup>1</sup> DTIS, ONERA, Université Paris Saclay F-91123 Palaiseau, France

<sup>2</sup> LIRMM, Univ. Montpellier, CNRS, Montpellier, France

## Corresponding author:

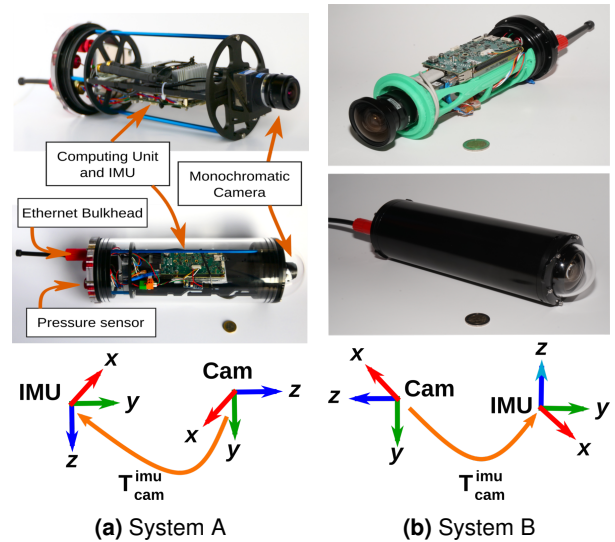
Maxime Ferrera

Email: maxime.ferrera@gmail.com

limits the development of new methods. [Bender et al. \(2013\)](#) proposed a dataset containing the measurements of navigational sensors, stereo cameras and a multibeam sonar. [Mallios et al. \(2017\)](#) released another dataset dedicated to localization and mapping in an underwater cave from sonar measurements. Images acquired by a monocular camera are also given for the detection of cones precisely placed in order to have a mean of estimating drift. However, in both datasets, the acquisition rate of the cameras is too low ( $<10$  Hz) for most of VSLAM and VIO methods. [Duarte et al. \(2016\)](#) created a synthetic dataset simulating the navigation of a vehicle in an underwater environment and containing monocular cameras measurements at a frame-rate of 10 Hz. Many public datasets have also been made available by the Oceanography community through national websites (<https://www.data.gov/>, <http://www.marine-geo.org>). However, these datasets have not been gathered with the aim of providing data suitable for VSLAM or VIO and often lack essential information such as the calibration of their sensors' setup.

In this paper, we present AQUALOC, a new dataset aiming at the development of VSLAM and VIO methods dedicated to the underwater environment. The dataset sequences have been recorded using acquisition systems composed of a monochromatic camera, a Micro Electro-Mechanical System (MEMS) based IMU, a pressure sensor and a computing unit for synchronous recordings. These acquisition systems have been embedded on ROVs equipped with lighting systems and navigating close to the seabed. The recorded video sequences exhibit the typical visual degradation induced by the underwater environment such as turbidity, backscattering, shadows and strong illumination shifts caused by the artificial lighting systems. Three different sites have been explored to create the dataset: a harbor and two archaeological sites. The recording of the sequences occurred at different depths, going from a few meters, for the harbor, to several hundred meters, for the archaeological sites. The provided video sequences are hence highly diversified in terms of scenes (low-textured areas, very texture repetitive areas...) and of scenarios (exploration, photogrammetric surveys, manipulations...). As the acquisition of ground truth is very difficult in natural underwater environments, we have used the state-of-the-art Structure-from-Motion (SfM) library Colmap ([Schönberger and Frahm \(2016\)](#)) to compute comparative baseline trajectories for each sequence. Colmap processes offline the sequences and performs a 3D reconstruction to estimate the positions of the camera. This 3D reconstruction is done by matching exhaustively all the images composing a sequence, which allows the detection of many loop-closures and, hence, the computation of accurate trajectories, assessed by low average reprojection errors. Along with the computed trajectories, we also provide the list of matched images for each sequence which could be used to evaluate relocalization or loop-closure detection methods. We further include statistics on the 3D reconstruction to assess their accuracy.

With the release of this dataset, we provide to the community the opportunity to work on data difficult to acquire. Indeed, the logistic (ship availability) and the required equipment (deep-sea compliant underwater vehicles



**Figure 1.** The acquisition systems equipped with a monocular monochromatic camera, a pressure sensor, an IMU and a computer along with the sensors' reference frames.

and sensors), as well as regulations (official authorizations), can be a barrier preventing possible works on this topic. We are convinced that the availability of this dataset will increase the development of algorithms dedicated to the underwater environment. Both raw and ROS bag formatted field data are provided along with the full calibration of the sensors (camera and IMU). Moreover, the provided comparative baseline makes this dataset suitable for benchmarking VSLAM and VIO algorithms.

The rest of this paper is organized as follows. First, we present the design of the acquisition systems used and the calibration procedures employed. Then, an overview of the dataset is given and the acquisition conditions on each site are detailed, highlighting the associated challenges for visual localization. Next, the processing of the data sequences to create a baseline is described. Finally, we detail how the dataset is organized and in which way the data are formatted.

## 2 The Acquisition Systems

In order to acquire the sequences of the dataset, we have designed two similar underwater systems. These acquisition systems have been designed to allow the localization of underwater vehicles from a minimal set of sensors in order to be as cheap and as versatile as possible. Both systems are equipped with a monochromatic camera, a pressure sensor, a low-cost MEMS-IMU and an embedded computer. The camera is placed behind an acrylic dome to minimize the distortion effects induced by the difference between water and air refractive indices. The image acquisition rate is 20 Hz. The IMU delivers measurements from a 3-axes accelerometer, 3-axes gyroscope and 3-axes magnetometer at 200 Hz. The embedded computer is a Jetson TX2 running Ubuntu 16.04 and is used to record synchronously the sensors' measurements thanks to the ROS middleware. The Jetson TX2 is equipped with a carrier board embedding the mentioned MEMS-IMU and a 1 To NVME SSD to directly store the sensors measurements, thus avoiding any

bandwidth or package loss issue. An advantage of the self-contained systems that we have developed, is that they are independent of any robotic architecture and can thus be embedded on any kind of Remotely Operated Vehicle (ROV) or Autonomous Underwater Vehicle (AUV). The interface can either be Ethernet or a serial link, depending on the host vehicle's features.

To record data at different depths, we have designed two systems that we will refer to as ‘‘System A’’ and ‘‘System B’’. These systems have the same overall architecture, but they differ on the camera model, the pressure sensor type and the diameter and material of the enclosure. System A (Fig. 1a) is designed for shallow waters and was used to acquire the sequences in the harbor. Its camera has been equipped with a wide-angle lens, which can be modelled by the fisheye distortion model. The pressure sensor is rated for 30 bars and delivers depth measurements at a maximum rate of 10 Hz. System A is protected by an acrylic enclosure, rated for a depth of 100 meters. System B (Fig. 1b) was used to record the sequences on the archaeological sites at larger depths. Its camera has a slightly lower field of view and the lens can be modelled by the radial-tangential distortion model. It embeds a pressure sensor rated for 100 bars delivering depth measurements at 60 Hz. Its enclosure is made of aluminum and is 400 meters depth rated. The technical details about both systems and their embedded sensors are given in table 1.

Each camera-IMU setup has been cautiously calibrated to provide the intrinsic and extrinsic parameters required to use it for localization purpose. We have used the toolbox Kalibr (Furgale et al. (2012, 2013)) along with an apriltag target to compute all the calibration parameters.

The cameras calibration step allows to obtain an estimate of the focal lengths, principal points and distortion coefficients. These parameters can then be used to undistort the captured images and to model the image formation pipeline, with the following notation:

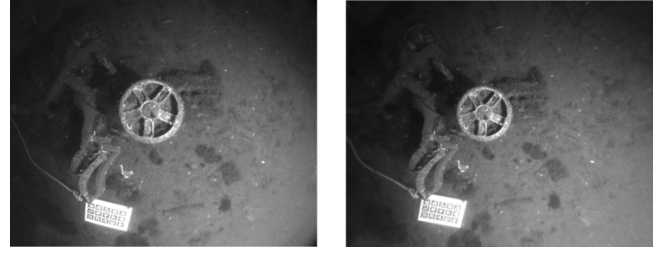
$$\begin{bmatrix} u \\ v \end{bmatrix} = \Pi_{\mathbf{K}} (\mathbf{R}_w^{\text{cam}} \mathbf{X}_w + \mathbf{t}_w^{\text{cam}}) \quad (1)$$

$$\begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} f_x \frac{x_{\text{cam}}}{z_{\text{cam}}} + c_x \\ f_y \frac{y_{\text{cam}}}{z_{\text{cam}}} + c_y \end{bmatrix} = \Pi_{\mathbf{K}} (\mathbf{X}_{\text{cam}}) \quad (2)$$

$$\text{with } \mathbf{K} = \begin{bmatrix} f_x & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix} \text{ and } \mathbf{X}_{\text{cam}} = \begin{bmatrix} x_{\text{cam}} \\ y_{\text{cam}} \\ z_{\text{cam}} \end{bmatrix}$$

where  $\Pi_{\mathbf{K}}(\cdot)$  denotes the projection:  $\mathbb{R}^3 \mapsto \mathbb{R}^2$ ,  $\mathbf{K}$  is the calibration matrix,  $\mathbf{X}_w \in \mathbb{R}^3$  is the position of a 3D landmark in the world frame,  $\mathbf{R}_w^{\text{cam}} \in SO(3)$  and  $\mathbf{t}_w^{\text{cam}} \in \mathbb{R}^3$  denote the rotational and translational components of the transformation from the world frame to the camera frame,  $\mathbf{X}_{\text{cam}} \in \mathbb{R}^3$  is the position of a 3D landmark in the camera frame,  $f_x$  and  $f_y$  denotes the focal lengths and  $(c_x, c_y)$  is the principal point of the camera.

As these parameters are medium dependant, the calibration has been performed in water to account for the additional distortion effects at the dome's level. The results of the calibration of the fisheye camera can be seen in figure 2.



**Figure 2.** Distortion effects removal from Kalibr calibration on one of the harbor sequences. Left: raw image. Right: undistorted image.

The camera-IMU setup calibration allows to estimate the extrinsic parameters of the setup, that is the relative position of the camera with respect to the IMU, and the time delay between camera's and IMU's measurements. This relative position is represented by a rotation matrix  $\mathbf{R}_{\text{cam}}^{\text{imu}}$  and a translation vector  $\mathbf{t}_{\text{cam}}^{\text{imu}}$ . Camera and IMU's poses relate to each other through:

$$\mathbf{T}_{\text{cam}}^w = \mathbf{T}_{\text{imu}}^w \mathbf{T}_{\text{cam}}^{\text{imu}} \quad (3)$$

$$\text{with } \mathbf{T}_{\text{cam}}^{\text{imu}} \doteq \begin{bmatrix} \mathbf{R}_{\text{cam}}^{\text{imu}} & \mathbf{t}_{\text{cam}}^{\text{imu}} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4}$$

$$\text{and } (\mathbf{T}_{\text{cam}}^w)^{-1} = \mathbf{T}_w^{\text{cam}} \doteq \begin{bmatrix} \mathbf{R}_w^{\text{cam}} & \mathbf{t}_w^{\text{cam}} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \in \mathbb{R}^{4 \times 4}$$

where  $\mathbf{R}_{\text{cam}}^{\text{imu}} \in SO(3)$ ,  $\mathbf{t}_{\text{cam}}^{\text{imu}} \in \mathbb{R}^3$ ,  $\mathbf{T}_{\text{cam}}^w \in SE(3)$ ,  $\mathbf{T}_w^{\text{cam}} \in SE(3)$ ,  $\mathbf{T}_{\text{cam}}^{\text{imu}} \in SE(3)$  and  $\mathbf{T}_{\text{imu}}^w \in SE(3)$ .  $\mathbf{T}_{\text{cam}}^w$  and  $\mathbf{T}_{\text{imu}}^w$  respectively represent the poses of the camera and of the body, with respect to the world frame.  $\mathbf{T}_w^{\text{cam}}$  is the inverse transformation of  $\mathbf{T}_{\text{cam}}^w$  and  $\mathbf{T}_{\text{cam}}^{\text{imu}}$  is the transformation from the camera frame to the IMU frame.

Before estimating these extrinsic parameters, the IMU noise model parameters have been derived from an Allan standard deviation plot, obtained by recording the gyroscope and accelerometer measurements for several hours, while keeping the IMU still. These noise parameters are then fed into the calibration algorithms to model the IMU measurements. As these parameters (IMU noises, camera-IMU relative transformation and measurements time delay) are independent of the medium (air or water), they have been estimated in air. Doing this calibration step in air allowed to perform easily the fast motions required to correlate the IMU measurements to the camera's ones.

All the calibration results are included in the dataset, that is the cameras' models (including the intrinsic parameters and the distortion coefficients), the IMUs' noise parameters, the relative transformation from the camera to the IMU and the time delay between the cameras' and the IMUs' measurements.

### 3 Dataset Overview

As explained in section 2, System A was used to record the shallow harbor sequences, while System B was used on the two deep archaeological sites. We propose a total

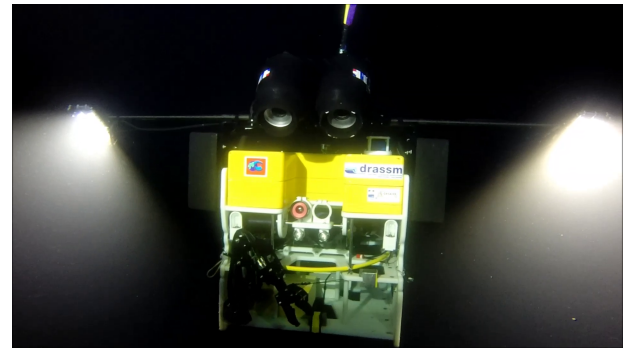
System A (Harbor sequences)	<b>Camera sensor</b>	<b>UEye - UI-1240SE</b>
	Resolution	640×512 px
	Sensor	Monochromatic
	Frames per second	20 fps
	<b>Lens</b>	<b>Kowa LM4NCL C-Mount</b>
	Focal length	3.5mm
	<b>Pressure Sensor</b>	<b>MS5837 - 30BA</b>
	Depth range	0 - 290m
	Resolution	0.2 mbar
	Output frequency	5-10 Hz
	<b>Inertial Measurement Unit</b>	<b>MEMS - MPU-9250</b>
	Gyroscope frequency	200 Hz
	Accelerometer frequency	200 Hz
	Magnetometer frequency	200 Hz
<b>Embedded Computer</b>	<b>Nvidia - Tegra Jetson TX2</b>	
Carrier board	Auvideo J120 - IMU	
Storage	NVME SSD 1 To	
<b>Housing</b>	<b>4" Blue Robotics Enclosure</b>	
Enclosure	33.4 x 11.4 cm	
Enclosure Material	Acrylic	
Dome	4" Blue Robotics Dome End Cap	
System B (Archaeo. sequences)	<b>Camera sensor</b>	<b>UEye - UI-3260CP</b>
	Resolution	968×608 px
	Sensor	Monochromatic
	Frames per second	20 fps
	<b>Lens</b>	<b>Kowa LM6NCH C-Mount</b>
	Focal length	6mm
	<b>Pressure Sensor</b>	<b>Keller 7LD - 100BA</b>
	Depth range	0 - 990m
	Resolution	3 mbar
	Output frequency	60 Hz
	<b>Inertial Measurement Unit</b>	<b>Same as System A</b>
	<b>Embedded Computer</b>	<b>Same as System A</b>
	<b>Housing</b>	<b>3" Blue Robotics Enclosure</b>
	Enclosure	25.8 x 8.9 cm
Enclosure Material	Aluminium	
Dome	3" Blue Robotics Dome End Cap	

**Table 1.** Technical details about the acquisition systems.



**Figure 3.** The Remotely Operated Vehicle *Dumbo* and the acquisition system A, used to record the harbor sequences.

of 17 sequences: 7 recorded in the harbor, 4 on the first archaeological site and 6 on the second site. As each of these environments is in some ways different from the others, we describe the sequences recorded in each environment separately. Table 2 summarizes the specificities of each data sequence. Note that, for each sequence, the starting and ending points are approximately the same. In most of the sequences, there are closed loops along the performed trajectories. Some sequences also slightly overlap, which can be useful for the development of relocalization features.



**Figure 4.** The Remotely Operated Vehicle *Perseo*, used on the archaeological sites.

*Credit: F. Osada - DRASSM / Images Explorations.*

### 3.1 Harbor sequences

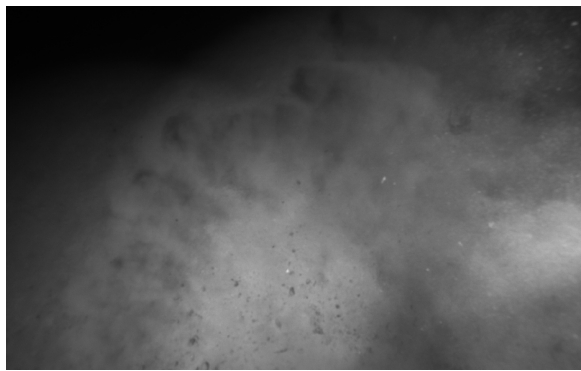
The harbor sequences were recorded in April 2018. System A was embedded on the lightweight ROV *Dumbo* (DRASSM-LIRMM) with the camera facing downward, as shown in figure 3. The ROV was navigating at a depth of 3 to 4 meters over an area of around 100 m<sup>2</sup>. Although the sun illuminates this shallow environment, a lighting system was used in order to increase the signal-to-noise ratio of the images acquired by the camera. The explored area was mostly planar but the presence of several big objects made it a real 3D environment, with significant relief.

For each sequence, loops are performed and an apriltag calibration target is used as a marker for starting and ending points. On these sequences, vision is mostly degraded by light absorption, strong illumination variations and backscattering. In two sequences, visual information even becomes unavailable for a few seconds because of collisions with surrounding objects. Another challenge is the presence of areas with seagrass moving because of the swell. Moreover, the ROV is sensitive to waves and tether disturbances, which results in roll and pitch variations.

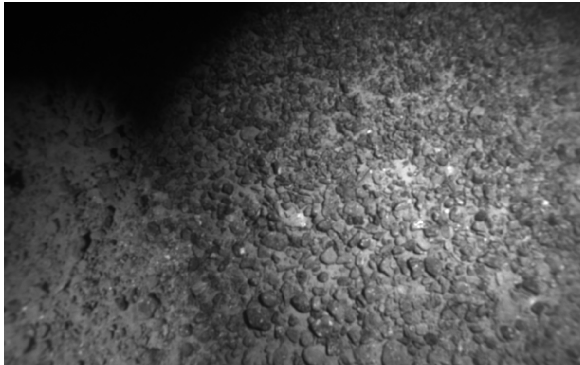
### 3.2 Archaeological sites sequences

The archaeological sites sequences were recorded in the Mediterranean sea, off Corsica's shore. The System B, designed for deep waters, was embedded on the *Perseo* ROV (Copetech SM Company) displayed in Fig. 4. In the way it was attached to the ROV, the camera viewing direction made a small angle with the vertical line ( $\approx 20 - 30^\circ$ ). *Perseo* is equipped with two powerful led lights (250,000 lumens each) and with two robotics arms for manipulation purposes. As localization while manipulating objects is a valuable information, to grab an artifact for instance, in some sequences the robotic arms are in the camera's field of view. A total of 10 sequences have been recorded on these sites, with 3 sequences taken on the first site and 7 on the second one.

The first archaeological site explored was located at a depth of approximately 270 meters and hosted the remains of an antic shipwreck. Hence, this site is mostly planar and presents mainly repetitive textures, due to numerous small rocks that were used as ballast in this antic ship (Fig. 5a). These sequences are affected by turbidity and moving sand particles, increasing backscattering and creating sandy

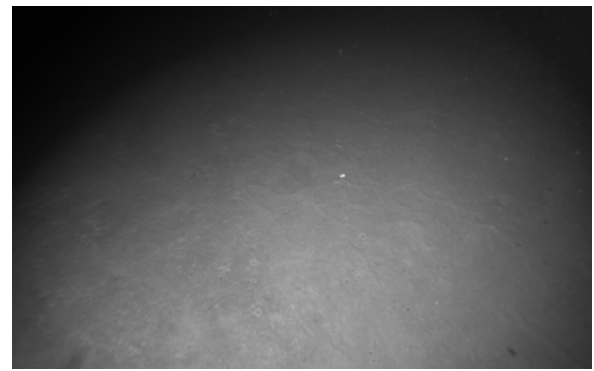


(a) Sandy cloud



(b) Texture repetitive area

**Figure 5.** Images acquired on the first archaeological site (depth: 270m).



(a) Low texture area



(b) Hill of amphorae

**Figure 6.** Images acquired on the second archaeological site (depth: 380m).

clouds (Fig. 5b). These floating particles are stirred up from the seabed by the water flows of the ROV's thrusters and lead to challenging visual conditions. A shadow is also omnipresent in these sequences in the left corner of the recorded images, because of the limits of the lighting system.

The second visited archaeological site was located at a depth of approximately 380 meters. On this site a hill of amphorae is present (Fig. 6b), whose top is culminating a few meters above the surrounding seabed level. During these sequences, the ROV was mainly operated for manipulation and photogrammetry purposes. While the amphorae present high texture, the ROV was also hovering low-textured sandy areas around the hill of amphorae (Fig. 6a). Because of the presence of these amphorae, marine wildlife has been growing on this site. Hence, the environment is quite dynamic, with many fishes getting in the field of view of the camera and many shrimps moving in the vicinity of the amphorae. In one of the sequence, both arms get in front of the camera. Otherwise, the visual degradation are the same as on the first site.

## 4 Comparative Baseline

As the acquisition of a ground truth is very difficult in natural underwater environment, we have used the state-of-the-art Structure-from-Motion (SfM) software Colmap (Schönberger and Frahm (2016)) to offline compute a 3D reconstruction for each sequence and extract a reliable trajectory from it. By setting very low the features extraction parameters, we were able to extract enough SIFT features (Lowe (2004)) to robustly match the images of each

sequence. Performing a matching of the images in an exhaustive way, that is trying to match each image to all the other ones, allows to get a reliable trajectory reconstruction as many closed loops can be found (Fig. 7). In Table 3, we provide statistics for each sequence about Colmap's 3D reconstructions to highlight the reliability of the reconstructed models. These statistics include the number of images used, the number of estimated 3D points, the average track length of each 3D points (*i.e.* the number of images observing a given 3D point) and the average reprojection error. The high average track lengths for each sequence (going from 6.7 to more than 20) assess the accuracy of the 3D points' estimation as it leads to a high redundancy in the bundle adjustment steps of the reconstruction. Moreover, given these high track lengths, the average reprojection error is a good indicator of the overall quality of a SfM 3D model and for each one of the sequences this error is below 0.9 pixel.

The extracted trajectories have been scaled using the pressure sensor measurements and hence provide metric positions. Although these trajectories cannot be considered as being perfect ground truths, we believe that it provides a fair baseline to evaluate and compare online localization methods. Evaluation of such methods can be done using the standard Relative Pose Error (RPE) and Absolute Trajectory Error (ATE) metrics (Sturm et al. (2012)).

Furthermore, we have made available the list of overlapping images (*i.e.* matching) according to Colmap for each sequence. These files could hence be used to evaluate the efficiency of loop-closure or image retrieval methods.

Site	Sequence	Duration	Length	Visual Disturbances					
				Turbidity	Collisions	Backscattering	Sandy clouds	Dynamics	Robotic Arm
Harbor (depth $\approx$ 4 m) Acquired by system A, embedded on a lightweight ROV	#01	3'49"	39.3m	X	-	X	-	-	-
	#02	6'47"	75.6m	X	-	X	-	-	-
	#03	4'17"	23.6m	X	-	X	-	-	-
	#04	3'26"	55.8m	X	X	X	-	-	-
	#05	2'52"	28.5m	X	-	X	-	-	-
	#06	2'06"	19.5m	X	-	X	-	-	-
	#07	1'53"	32.9m	X	X	X	-	-	-
First Archaeological Site (depth $\approx$ 270 m) Acquired by System B, embedded on a medium workclass ROV	#01	14'39"	32.4m	X	-	X	X	X	X
	#02	7'29"	64.3m	X	-	X	X	X	-
	#03	5'16"	10.7m	X	-	X	X	-	-
Second Archaeological Site (depth $\approx$ 380 m) Acquired by System B embedded on a medium workclass ROV	#04	11'09"	18.1m	X	-	X	X	X	X
	#05	3'19"	42.0m	X	-	X	-	X	-
	#06	2'49"	31.8m	X	-	X	-	X	-
	#07	9'29"	122.1m	X	-	X	-	X	-
	#08	7'49"	41.2m	X	-	X	-	X	-
	#09	5'49"	65.4m	X	-	X	-	X	-
	#10	11'54"	83.5m	X	-	X	-	X	-

**Table 2.** Details on all the AQUALOC sequences and their associated visual disturbances.

## 5 Data Sequences Format

As explained in the introduction, the sequences are all available as ROS bags and as raw data. The dataset is split into two folders, one for the harbor sequences and the other for the archaeological ones.

The dataset repository architecture is the following:

```

Harbor_site_sequences/
├── Calibration files/
│   ├── camera_calib.txt
│   ├── imu_camera_calib.txt
│   └── imu_noises.txt
├── ground truth files/
│   ├── colmap_traj_sequence_X.txt
│   ├── ...
│   ├── colmap_detected_loop_sequence_X.txt
│   └── ...
├── sequence_X_bag.tar.gz/
│   └── sequence_X.bag
│   └── ...
├── sequence_X_raw_data.tar.gz/
│   ├── imu.csv
│   ├── mag.csv
│   ├── images.csv
│   └── images/
│       └── frameXXXXXX.png
│       └── ...

```

The archaeological sites sequences do not appear here but are organized exactly in the same manner.

The calibration files are given in the output format of Kalibr (Furgale et al. (2012, 2013)).

The trajectories computed by Colmap for each sequence are available as text files and contain the pose in a translation-quaternion form. These files format is the following:

#Frame	tx	ty	tz	qx	qy	qz	qw
0.	-1.88	2.41	-0.47	0.01	0.06	0.14	0.91
20.	-1.83	2.35	-0.46	0.05	0.64	0.14	0.99
40.	-1.80	2.10	-0.34	0.04	0.58	0.12	0.98
...							

The files containing the loop closures detected by Colmap provide information in the following format:

```

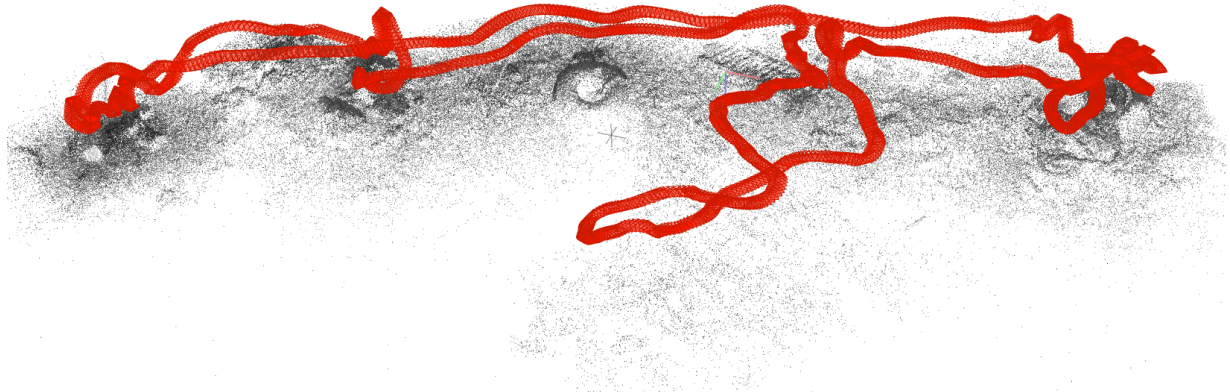
1, 1, 0, 0, 1
1, 1, 1, 0, 0
0, 1, 1, 0, 0
0, 0, 0, 1, 1
1, 0, 0, 1, 1
...

```

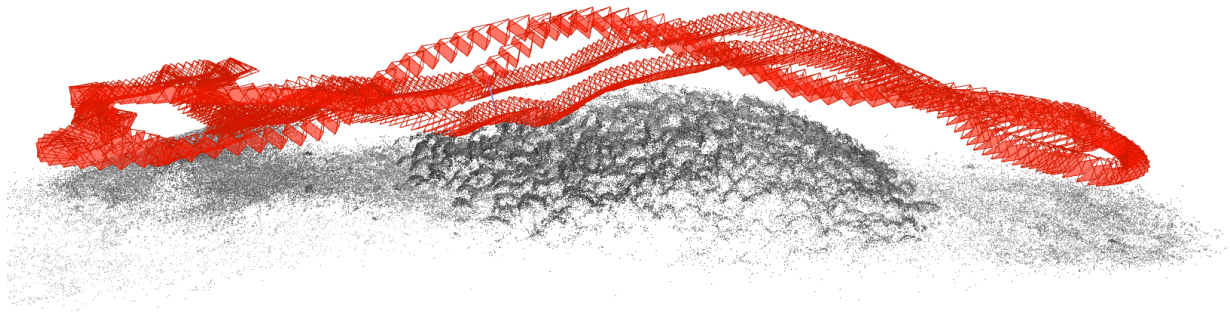
where a 1 indicates an overlapping between row  $i$  and column  $j$ , with  $i$  and  $j$  standing for the frame numbers. Note that only a subset of the images has been used to compute the offline reconstruction with Colmap (1 image out of 5 for the harbor sequences and 1 out of 20 for the archaeological ones). Therefore, the frame number given in these ground truth files is the number of their corresponding frame in the full sequence.

About the bag files, each sequence is stored in a separate bag containing the following topics:

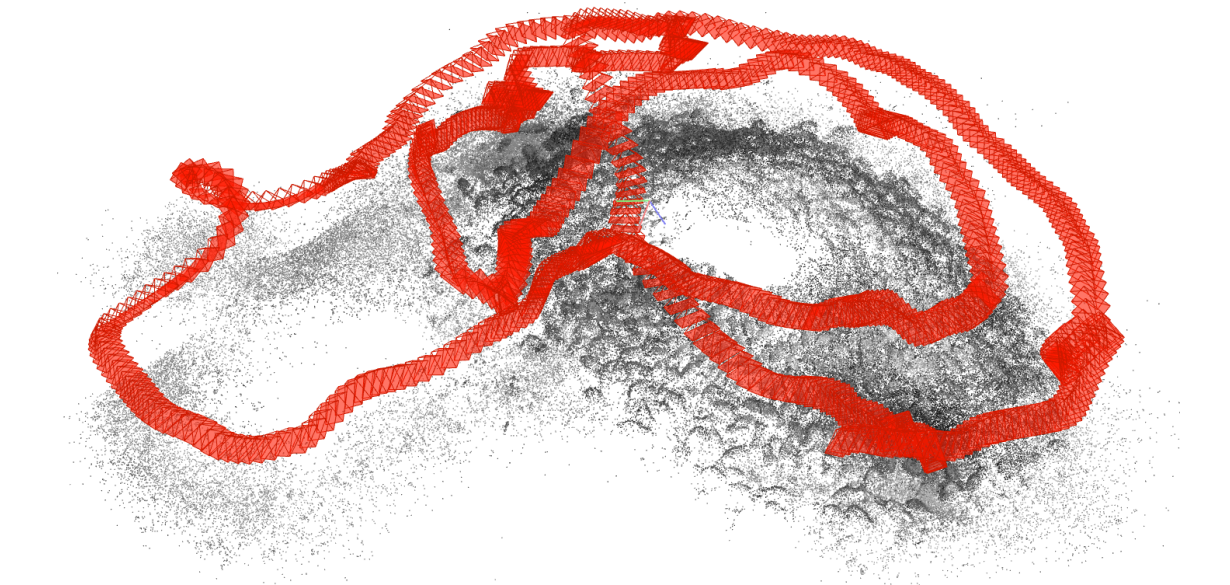
- **/camera/image\_raw**: Images recorded from the camera.
- **/camera/camera\_info**: Images width and height info.
- **/rtimulib\_node/imu**: Accelerometer and gyroscope measurements.
- **/rtimulib\_node/mag**: Magnetometer measurements.
- **/barometer\_node/pressure**: Pressure measurements in millibars.
- **/barometer\_node/depth**: Depth measurements in meters.



(a) Colmap reconstruction - Harbor #02



(b) Colmap reconstruction - Archeological Site #07



(c) Colmap reconstruction - Archeological Site #10

Figure 7. Examples of trajectories reconstructed with Colmap.

	Harbor sequences							Archeological sites sequences									
	#01	#02	#03	#04	#05	#06	#07	#01	#02	#03	#04	#05	#06	#07	#08	#09	#10
Nb. of used images	918	1590	1031	770	692	508	447	880	445	311	637	200	170	569	470	350	715
Nb. of 3D points	112659	305783	355130	194407	236845	188807	181964	196857	174514	160531	249048	42877	45799	251620	237882	114814	329686
Mean tracking length	14.9	13.2	17.2	9.7	10.7	12.1	9.5	23.5	12.6	8.4	8.5	7.6	6.7	7.4	9.1	7.9	9.2
Mean reproj. err. (px)	0.896	0.816	0.713	0.715	0.688	0.733	0.846	0.746	0.621	0.474	0.673	0.601	0.569	0.645	0.616	0.660	0.661

Table 3. Colmap trajectories reconstruction statistics. The number of provided images, the number of reconstructed 3D points, the mean tracking length for the 3D points and the mean reprojection error for the 3D reconstruction are given for each sequence.



- **/barometer\_node/temperature**: Pressure sensor's temperature measurements.

In their raw format, each sequence contains the following data:

- **:** The directory containing the sequence images.
- **frameXXXXX.png**: The images recorded from the camera.
- **images.csv**: The timestamps related to each image of the sequence.
- **imu.csv**: The accelerometer and gyroscope measurements and their timestamps.
- **mag.csv**: The magnetometer measurements and their timestamps.
- **depth.csv**: The pressure measurements converted in meters and their timestamps.

For each *csv* files, the first row starts with a # and then gives the name of the different fields along with their related measurements unit into squared brackets. The following rows contain the values of the measurements. In all these files, the first field is the acquisition timestamp of the measurements. For instance, the *depth.csv* files look like:

```
#timestamp [ns], depth [m]
1542828791719540119,271.988866935
1542828791735507011,272.01910918
...
```

## 6 Conclusion

In this paper, we have presented a new dataset of subsea monocular video sequences synchronized with inertial and pressure measurements. This dataset is intended for encouraging the development of localization methods for underwater robots navigating close to the seabed. The sequences have been recorded from Remotely Operated Vehicles in three different environments at different depths: a harbor at a depth of 4 meters, a first archaeological site at a depth of 270 meters and a second one at a depth of 380 meters. The diversity of the recorded environments allowed to capture video sequences with different visual perturbations typical in underwater scenarios. For each sequence, trajectories have been computed offline using a Structure-from-Motion library and are provided as a baseline for performance comparisons of localization methods. The datasets are available both as ROS bags and as raw data. In future work, we plan to perform new acquisition missions in different underwater environments in order to augment this dataset and increase its diversity.

## Acknowledgements

The experiments conducted on the archaeological sites have been done in the French waters, in the framework of a research campaign of the DRASSM (French Department of Underwater Archaeology - Ministry of Culture) in accordance with the international regulation (UNESCO Convention on the Protection of the Underwater Cultural Heritage, 2001).

The authors are grateful to the DRASSM for its logistical support. The authors acknowledge support of the CNRS (Mission pour l'interdisciplinarité - Instrumentation aux limites 2018 - Aqualoc project) and support of Région Occitanie (ARPE Pilotplus project).

The authors would also like to thank Anthelme Bernard-Brunel for his help in the design of the acquisition systems and Abderrahmane Kheddar for his helpful remarks.

## 6.1 References

### References

- Bender A, Williams SB and Pizarro O (2013) **Autonomous exploration of large-scale benthic environments**. In: *2013 IEEE International Conference on Robotics and Automation (ICRA)*. Karlsruhe, Germany, pp. 390–396.
- Blanco JL, Moreno FA and Gonzalez-Jimenez J (2014) **The Málaga Urban Dataset: High-rate Stereo and Lidars in a realistic urban scenario**. *International Journal of Robotics Research* 33(2): 207–214.
- Burguera A, Bonin-Font F and Oliver G (2015) **Trajectory-Based Visual Localization in Underwater Surveying Missions**. *Sensors* 15(1): 1708–1735.
- Burri M, Nikolic J, Gohl P, Schneider T, Rehder J, Omari S, Achtelik MW and Siegwart R (2016) **The EuRoC micro aerial vehicle datasets**. *The International Journal of Robotics Research* 35(10): 1157–1163.
- Cadena C, Carlone L, Carrillo H, Latif Y, Scaramuzza D, Neira J, Reid I and Leonard JJ (2016) **Past, Present, and Future of Simultaneous Localization and Mapping: Toward the Robust-Perception Age**. *IEEE Transactions on Robotics* 32(6): 1309–1332.
- Creuze V (2017) **Monocular Odometry for Underwater Vehicles with Online Estimation of the Scale Factor**. In: *IFAC 2017 World Congress*. Toulouse, France.
- Duarte AC, Zaffari GB, da Rosa RTS, Longaray LM, Drews P and Botelho SSC (2016) **Towards comparison of underwater SLAM methods: An open dataset collection**. In: *OCEANS 2016 MTS/IEEE*. Monterey, CA, USA, pp. 1–5.
- Engel J, Koltun V and Cremers D (2018) **Direct Sparse Odometry**. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 40(3): 611–625.
- Ferrera M, Moras J, Trouvé-Peloux P and Creuze V (2019) **Real-Time Monocular Visual Odometry for Turbid and Dynamic Underwater Environments**. *Sensors* 19(3).
- Forster C, Zhang Z, Gassner M, Werlberger M and Scaramuzza D (2017) **SVO: Semidirect Visual Odometry for Monocular and Multicamera Systems**. *IEEE Transactions on Robotics* 33(2): 249–265.
- Furgale P, Barfoot T and Sibley G (2012) **Continuous-Time Batch Estimation Using Temporal Basis Functions**. In: *2012 IEEE International Conference on Robotics and Automation (ICRA)*. St. Paul, MN, USA, pp. 2088–2095.
- Furgale P, Rehder J and Siegwart R (2013) **Unified Temporal and Spatial Calibration for Multi-Sensor Systems**. In: *2013 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Tokyo, Japan, pp. 1280–1286.
- Geiger A, Lenz P and Urtasun R (2012) **Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite**. In: *2012 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Providence, RI, USA, pp. 3354–3361.
- Kim A and Eustice RM (2013) **Real-Time Visual SLAM for Autonomous Underwater Hull Inspection Using Visual Saliency**. *IEEE Transactions on Robotics* 29(3): 719–733.

- Leutenegger S, Lynen S, Bosse M, Siegwart R and Furgale P (2015) **Keyframe-based visual-inertial odometry using nonlinear optimization**. *The International Journal of Robotics Research* 34(3): 314–334.
- Lowe DG (2004) **Distinctive Image Features from Scale-Invariant Keypoints**. *International Journal of Computer Vision* 60(2): 91–110.
- Mallios A, Vidal E, Campos R and Carreras M (2017) **Underwater caves sonar data set**. *The International Journal of Robotics Research* 36(12): 1247–1251.
- Mur-Artal R, Montiel JMM and Tardós JD (2015) **ORB-SLAM: A Versatile and Accurate Monocular SLAM System**. *IEEE Transactions on Robotics* 31(5): 1147–1163.
- Mur-Artal R and Tardos JD (2017) **Visual-Inertial Monocular SLAM With Map Reuse**. *IEEE Robotics and Automation Letters* 2(2): 796–803.
- Paull L, Saeedi S, Seto M and Li H (2014) **AUV Navigation and Localization: A Review**. *IEEE Journal of Oceanic Engineering* 39(1): 131–149.
- Qin T, Li P and Shen S (2018) **VINS-Mono: A Robust and Versatile Monocular Visual-Inertial State Estimator**. *IEEE Transactions on Robotics* 34(4): 1004–1020.
- Quattrini Li A, Coskun A, Doherty SM, Ghasemlou S, Jagtap AS, Modasshir M, Rahman S, Singh A, Xanthidis M, O’Kane JM and Rekleitis I (2017) **Experimental Comparison of Open Source Vision-Based State Estimation Algorithms**. In: *2016 International Symposium on Experimental Robotics*. pp. 775–786.
- Rahman S, Li AQ and Rekleitis I (2018) **Sonar Visual Inertial SLAM of Underwater Structures**. In: *2018 IEEE International Conference on Robotics and Automation (ICRA)*. Brisbane, QLD, Australia, pp. 1–7.
- Schönberger JL and Frahm JM (2016) **Structure-from-Motion Revisited**. In: *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Las Vegas, NV, USA.
- Shkurti F, Rekleitis I, Scaccia M and Dudek G (2011) **State estimation of an underwater robot using visual and inertial information**. In: *2011 IEEE/RSJ Intelligent Robots and Systems (IROS)*. San Francisco, CA, USA.
- Sturm J, Engelhard N, Endres F, Burgard W and Cremers D (2012) **A benchmark for the evaluation of RGB-D SLAM systems**. In: *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Vilamoura, Portugal, pp. 573–580.
- Weidner N, Rahman S, Li AQ and Rekleitis I (2017) **Underwater cave mapping using stereo vision**. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. Singapore, Singapore, pp. 5709–5715.
- Zhang J, Ila V and Kneip L (2018) **Robust Visual Odometry in Underwater Environment**. In: *2018 OCEANS MTS/IEEE Kobe Techno-Oceans (OTO)*. Kobe, Japan, pp. 1–9.