

Received February 15, 2022, accepted March 8, 2022, date of publication March 11, 2022, date of current version March 21, 2022.

Digital Object Identifier 10.1109/ACCESS.2022.3158971

DOES: A Deep Learning-Based Approach to Estimate Roll and Pitch at Sea

FABIANA DI CIACCIO¹, PAOLO RUSSO², AND SALVATORE TROISI¹

¹International Ph.D. Program/UNESCO Chair "Environment, Resources and Sustainable Development," Department of Science and Technology, Parthenope University of Naples, 80143 Naples, Italy

²Department of Computer, Control and Management Engineering Antonio Ruberti, University of Rome La Sapienza, 00185 Rome, Italy

Corresponding author: Fabiana Di Ciaccio (fabiana.diciaccio@studenti.uniparthenope.it)

ABSTRACT The use of Attitude and Heading Reference Systems (AHRS) for orientation estimation is now common practice in a wide range of applications, e.g., robotics and human motion tracking, aerial vehicles and aerospace, gaming and virtual reality, indoor pedestrian navigation and maritime navigation. The integration of the high-rate measurements can provide very accurate estimates, but these can suffer from errors accumulation due to the sensors drift over longer time scales. To overcome this issue, inertial sensors are typically combined with additional sensors and techniques. As an example, camera-based solutions have drawn a large attention by the community, thanks to their low-costs and easy hardware setup; moreover, impressive results have been demonstrated in the context of Deep Learning. This work presents the preliminary results obtained by DOES, a supportive Deep Learning method specifically designed for maritime navigation, which aims at improving the roll and pitch estimations obtained by common AHRS. DOES recovers these estimations through the analysis of the frames acquired by a low-cost camera pointing the horizon at sea. The training has been performed on the novel ROPIS dataset, presented in the context of this work, acquired using the FrameWO application developed for the scope. Promising results encourage to test other network backbones and to further expand the dataset, improving the accuracy of the results and the range of applications of the method as a valid support to visual-based odometry techniques.

INDEX TERMS AHRS, computer vision, dataset acquisition, deep learning, orientation estimation.

I. INTRODUCTION

The pose estimation problem consists in estimating the position and orientation of a vehicle, device, human or robot with respect to a reference frame, through the use of different kinds of internal or external sensors. The accurate measurement of the orientation plays in fact a critical role in a wide range of activities, e.g., robotics and human motion tracking, bio-logging for animal behaviour research, aerial vehicles and aerospace, gaming and virtual reality applications, medicine and biotechnology, indoor and outdoor pedestrian navigation, maritime and/or autonomous navigation. When Global Navigation Satellite Systems (GNSS) are not able to provide correct information about the position and attitude of a vehicle, navigation and localization operations are generally performed through the integration of different kind of sensors: inertial, odometry, laser and sonar ranging sensors, underwater positioning systems, etc. [1].

The associate editor coordinating the review of this manuscript and approving it for publication was Pinjia Zhang.

In the last years the use of low-cost technologies is becoming widely spread in numerous applications: this means that the accuracy of the pose obtained by these systems can be affected by even more disturbing factors than the traditional high-performing methods. In these circumstances, the development of accurate and reliable orientation estimation algorithms can still be considered a very challenging task, being at the basis of the localization process and of the consequent performances of the device employed for any specific task. This finds particular application in the context of the navigation, be it aerial, maritime or pedestrian, underwater/underground or in surface, autonomous, remotely operated or traditionally performed. In the specific case of maritime navigation, the information of position and orientation of a vessel is of great interest for seafarers in different operations and scenarios (e.g., open sea, congested harbours and waterways) as it is strictly related to the safety of the navigation at any level [2]. The same goes for Unmanned Surface Vehicles (USVs), which are mainly employed in environmental monitoring, safety or navigation support and

research operations. In this case, a non accurate estimation of the orientation can severely compromise the ultimate success of the mission, especially when paired to low-cost sensors and poor GNSS support.

The Inertial Measurement Unit (IMU) gives the instantaneous speed and position of the vehicle without the need for external references by integrating the measures of angular velocity and linear acceleration obtained through its three orthogonal rate-gyroscopes and $-$ accelerometers respectively. Unfortunately, several problems are associated with these sensors; among the others, measurements are noisy and biased and the errors increase over time due to the drift of the sensors. Micro Electro-Mechanical Systems (MEMS) Attitude Heading Reference Systems (AHRS) integrate to this configuration a magnetometer which measures the variation of the Earth's magnetic field: this allows to instantly calculate an improved estimation whereas benefitting from lighter weight, smaller sizes and lower prices. The great potential of these devices makes them suitable for several applications exploiting the pure orientation estimation, like geomatics, surveys, augmented reality, etc.

Vision-based methods are also frequently employed for the scope: these techniques allow to understand the surrounding environment by detecting its visual features through a camera; captured color data with their high resolution contain in fact several information, and the sensors are generally low-costs and with an easy hardware setup. In this context, the detection of the horizon line is an important attribute for the maritime image processing, as it allows to estimate the camera's orientation with respect to the sea surface other than restricting the object search region when detection is performed, thus reducing the processing time and the false detection problem. Several approaches have been proposed to solve this task, however the accuracy and the processing time of the horizon line detection on high-resolution maritime image still face some issues [3].

In the last decade, Visual Odometry (VO) and Visual Simultaneous Localization and Mapping (VSLAM) techniques have been successfully developed; however, their application can be challenging too, especially when their deployment is made in non-textured environments or with poor-light conditions. To reduce these limitations, IMU and camera systems are integrated in Visual Inertial Odometry (VIO) techniques [4]; as a drawback, they require manual interference for possible failure cases assessment, careful and specific tuning of the parameters related to the environment, and a final refining of the results. In recent years, increasing consideration has been gained by Deep Learning (DL) techniques, which demonstrated to be robust to camera parameters and harsh scenarios: these methods are in fact able to successfully extrapolate and learn new features representations from the images they are fed with and these can further improve the motion estimation [5].

With the aim of providing further enhancements in the orientation estimation methodologies, this paper presents DOES, Deep Orientation (of roll and pitch) Estimation at



FIGURE 1. Illustration of an image from the ROLL and Pitch at Sea (ROPIS) dataset.

Sea, a new supportive DL model which can be combined to the actual low-cost IMU-based configuration. This approach is not intended to substitute the current systems, but aims at improving the robustness of traditional methods when some limitations occur: the unavailability of GPS signals in indoor and under-surface environment, the undesirable high drift of inertial sensors in case of extended GPS outages and the issues of possible confusion with nearby robots for SONAR & RADAR are some of the limitations associated with these navigation systems. Visual-based methods help in this sense, since they constitute a powerful tool to estimate the pose of a camera through which the motion information is further recovered. These techniques can be classified as geometric or learning based: in the first case the camera geometry is explored to estimate the motion, whereas in the latter the model is fed with labeled data and then trained to accomplish the same task. The advantage of the learning-based methods is that they do not require the knowledge of the camera parameters and can estimate the orientation with correct scale even for monocular cases [6]. Moreover, visual methods can be further integrated with traditional, IMU-based orientation estimation algorithms to obtain a robust and reliable visual-inertial odometry system [7]. The work presented in this paper develops an affordable visual, learning-based backbone which estimates the attitude of a monocular camera which will be mounted on a vehicle.

The idea behind DOES is in fact to train a DL model able to output the vehicle attitude (in terms of roll and pitch angles) by processing the sea horizon view recorded by a low-cost camera. In particular, the latter needs to be mounted on the surface of an autonomous robot (or, similarly, on the bridge of traditional ships) with its axis parallel to the vehicle longitudinal axis, to correctly frame the horizon line. A similar approach could be further tested on Unmanned Aerial Vehicles (UAVs) too. To lay the foundation for this task, preliminary intensive tests have been conducted to verify the validity of the approach. Different DL architectures have been tested for the processing of the images acquired through an Android smartphone's camera.

In this context, the lack of datasets specifically designed for DL-based orientation estimation at sea has been evidenced. While tackling this issue, the need of acquisition methods assuring the synchronism of the measurements for a reliable

Ground Truth (GT) has been addressed too. For this reason, this paper presents also the first release of the Roll and Pitch at Sea (ROPIS) dataset (Fig. 1), which has been created through FrameWO, an Android application developed for the scope. The choice of employing low-cost sensors meets the necessity to develop affordable and smart tools to enhance the orientation estimation; for this reason, the first deployment of the dataset has been acquired using open-source libraries and software. In this preliminary release, the operating user acquires the data in the proximity of the seashore trying to simulate the real behaviour of a ship in navigation.

The aim of this project is to provide a supportive visual-based, low-cost technique for attitude estimation which can be easily deployed in the context of navigation at sea or other challenging scenarios, as it does not need to take into account camera models or related calibration issues.

More in detail, the main contributions of this work can be summarized as follows:

- The development of FrameWO, an Android smartphone application for the simultaneous acquisition of camera images and their corresponding device orientation.
- The release of ROPIS dataset, consisting of 22173 RGB images/Euler angles samples acquired with FrameWO application on eight different sea locations.
- A Deep Learning-based method to perform attitude estimation using horizon-depicting frames; DOES is specifically trained on the ROPIS dataset and provides fast and reliable estimations, further encouraging to operate for its deployment in real-time scenarios.

The paper is organized as follows: Section II gives a brief overview on the existing literature on the orientation estimation task exploited through different traditional, visual and DL-based methods; Section III gives a theoretical foundation to the subject, introducing the attitude estimation problem to further describe the DL architectures which best fit the task. In Section IV the ROPIS dataset will be presented, highlighting the issues and solutions encountered during the app creation and the data acquisitions. Section V details the experiments performed on DOES whereas the obtained results will be presented and discussed in Section VI; final considerations and future objectives will conclude the work in Section VII.

II. RELATED WORKS

The accurate measurement of the orientation plays a critical role in a wide range of activities. AHRS sensors (i.e. accelerometers, gyroscopes and magnetometers) provide reliable measurements whose integration gives accurate information about the pose (position and attitude) of any object they are rigidly attached to. In the last decade, traditional methods have seen a huge improvement due to the integration with different kind of sensors, aiming at reducing the inertial-related error accumulation and the costs whilst enhancing the robustness of the methodology. As previously mentioned, one of the most effective integration is made

through visual-based method, leveraging the potential of visual features and the low-cost of the devices. The following paragraphs give a concise review of the existing literature in the field of orientation estimation.

A. INERTIAL-BASED METHODS

There exists a large amount of literature on the use of inertial sensors for position and orientation estimation. The reason for this is related to their robust algorithms and their accurate solutions which makes them suitable for being used in several fields. Interestingly, relatively simple position and orientation estimation algorithms work quite well in practice, even if the model choice can sensibly affect the accuracy of the estimates [8].

There is a large and ever-growing number of application areas for inertial sensors, as for example robotics and human motion tracking [9], [10], bio-logging for animal behavior research [11], aerial vehicles and aerospace [12], [13], gaming, virtual reality and indoor pedestrian navigation [14]–[16], etc. In fact, the use of accurate inertial sensors and magnetic compasses was first introduced in the navigation field, but along with the development of MEMS technology, low-cost and small-size inertial and magnetic compass sensors appeared in various kinds of consumer electronics, game consoles, virtual reality applications and so on. The orientation representations and sensor fusion still remain the challenges to overcome [17]. Real-time orientation estimation algorithms based on low-cost IMU are analyzed in [18], where the approach is based on the relationships between the quaternion representing the platform orientation and the measurements of the sensors and the integration is performed through an Extended Kalman Filter (EKF). Researchers in [19] developed a low-cost and low-weight attitude estimator for autonomous helicopters based on an inclinometer and a gyroscope, while fusing the data coming from the sensors through a classic complementary filter; in [20] a gyro-free, quaternion-based attitude determination system which exploits low cost sensors is presented. Reference [21] implemented a complementary filter able to infer Micro Aerial Vehicle (MAV) attitude from observations of gravity and magnetic field, with the final algorithm able to work with both IMU and MARG sensors. Authors in [13] exploited an AHRS device together with a Unscented Kalman Filter algorithm to perform attitude estimation on UAVs. The same filter has been used in [22], which developed a novel navigation system for autonomous underwater vehicles that works without the presence of a GPS device, not available in underwater scenarios. Researchers in [23] proposed an Adaptive Kalman Filter which is able to provide pose estimations based on low-cost AHRS devices, whereas [24] and [25] investigated the use of AHRS in smartphones as cheap but reliable devices for angles estimation. A novel error-state Kalman filter is presented in [26], which provides highly accurate IMU orientation estimates which result to be robust to fluctuations in the registered local magnetic field or caused by abrupt movements. An indoor pedestrian

navigation method based on shoe-mounted MEMS IMU and ultra-wideband is discussed in [27], which used a quaternion-based Kalman Filter to integrate the data and to reduce the complexity of the method. In [28] a new orientation estimation strategy for a non-accelerated platform is presented. Based on a low-cost IMU, this method sees a nonlinear Luenberger observer estimating the angles and a recursive least-square algorithm calibrating the common magnetometer offsets. Authors in [29] describes a calibration method for MEMS IMU mounted on electric bicycles that can be made in real-time thanks to its independence to sensor biases and its a very low computation cost.

B. VISION-BASED METHODS

The possibility to employ visual data to perform orientation and in general pose estimation has been widely deepened in the past decades. Many researches have been focused on the horizon line detection, due to its relevance for visual geo-localization, port security, etc. However, some special features in real marine environments (e.g., clouds clutter, sea glint and weather conditions) frequently result in different kinds of interference in optical images. Authors in [30] proposed a Sea-Sky Line (SSL) detection method for USVs based on the computation of the gradient saliency, through which the line features of the SSL are effectively enhanced while other disturbances are attenuated. The SSL identification is achieved according to regions contrast, line segment length and orientation features, and optimal state estimation of SSL detection is implemented by a cubature Kalman filter. In [31] a fast method for detecting the horizon line in maritime scenarios is presented. It combines a multi-scale approach and a region-of-interest (ROI) detection, which allows to efficiently reduce the required processing information amount. A single edge map is then produced and the Hough transform and a least-square method are sequentially applied to accurately estimate the horizon line. The Hough transform is also used in [32], which proposed a sea-sky line detection system based on the local Otsu segmentation; similarly, authors in [33] recognized the horizon line in maritime images through a two-phase, coarse-fine detection algorithm which increases the overall method robustness. Another quick horizon line detection method is proposed in [34], which extracts the horizon line in real maritime image with improved reliability and faster execution with respect to other competitors. The horizon detection through vision sensors is also frequently exploited to obtain redundant orientation information in the field of unmanned aerial navigation. For example, authors in [35] proposed two attitude estimation methods: the first one searches for the best line fitting the horizon in thermal images, which allows to further estimate the pitch and roll angles using an infinite horizon line model. The second method exploits a Convolutional Neural Network (CNN) which predicts the angles on the basis of the raw pixel intensities from the same kind of images.

However, these methods alone cannot be considered totally robust and reliable, since the position and slope of the horizon are strictly related to the camera intrinsic (i.e., focal length, optical center, pixel aspect ratio and skew) and extrinsic (rotation and translation) parameters and to the model used to parametrize them. In [36] the authors surveyed a plethora of methods which perform pose estimation by fusing visual, inertial and magnetic measurements, integrating them through the use of an EKF. The combined use of IMU and vision information has been explored by [1], which exploits SURF visual features together with accelerometer and gyroscope data to retrieve the robot pose in an indoor setting. A comprehensive analysis of the behaviour of these features when used for visual odometry can be found in [37].

VO, VIO and SLAM algorithms have recently received much attention for their efficient and accurate ego-motion estimation in robotics. A VIO algorithm for the estimation of the motion state of UAVs with high accuracy is presented in [38]. Visual data and pre-integrated inertial measurements are here integrated in an optimization framework; the stable initialization of scale and gravity through pose constraints together with a local scale parameter allowed to take into account the uncertainty of the VIO initialization.

The use of stereo camera sensors for VO is a reliable and low-cost way for attitude estimation, but may encounter problems when deployed underwater. This setting is in fact characterized by poor imaging and usually inconsistent motion due to the water flow. This issue has been tackled by [39], which proposed an AUV localization technique based on a stereo underwater VO system to overcome the aforementioned difficulties. In the context of underwater robotics, [40] presented another VO method which demonstrated to be robust to visual perturbations in many challenging scenarios. In [41] a novel key-frame based SLAM system is proposed, where a robust initialization aims at refining the scale through the use of depth measurements. Together with an improved image quality and a fast preprocessing step, this demonstrated to solve the localization drift and loss issues. A monocular VI-SLAM algorithm providing accurate and robust motion tracking is presented in [42]. This is developed in two parallel thread: the first one deals with the EKF motion tracking updated through a consistent map to reduce the drift. In the second one, a visual-inertial bundle adjustment is performed on the obtained global maps to optimize the overall results. ORB-SLAM3 [43] is another worth mentioning method in this context. It allows to use both stereo and monocular RGB-D cameras in the VI and SLAM approach, ensuring a robust real-time operativity in any kind of environment thanks to the Maximum-a-Posteriori estimation.

The rise of Deep Learning, with powerful architectures able to tackle complex tasks such as classification [44], detection [45], segmentation [46], denoising [47], super resolution [48], has definitely changed the way vision data are exploited for pose estimation. Instead of relying on engineered, fixed features (e.g. SIFT [49], SURF [50]),

recent algorithms exploit deep networks as powerful features extractors or by directly estimating the pose vector in an end-to-end model, from input images to the output prediction. For example, in order to estimate camera orientation, [51] exploited a LSTM deep network together with a linear Kalman Filter to combine IMU and camera data, whereas in DeepVIO [5] the authors fused 2D optical flow features together with standard inertial data, obtaining state of the art results on KITTI [52] and EuRoC [53] datasets. The combination of a traditional IMU with a LIDAR laser scan has been proposed in [54], where a recurrent CNN perform this aggregation on a scan-to-scan basis. In [55] researchers proposed a method to estimate a camera six degrees of freedom and absolute scale by exploiting unsupervised data, getting good results in terms of pose accuracy on KITTI benchmark. In [56], the authors developed a generative framework able to exploit a GAN [57] model on unlabelled RGB images for 6-DoF pose camera motion prediction, demonstrating the efficacy of their approach both on KITTI and Cityscapes [58] datasets. The former method has been improved in [59] with a stack of GAN layers which demonstrated to be effective on ego-motion estimation tasks. A comprehensive review of the state of the art deep models for pose estimation can be found in [60].

III. METHOD

This section aims at providing a theoretical background to fully understand the fundamentals of the proposed work. In particular, a general overview on the orientation estimation process is given in subsection III-A, with some details on the sensors embedded in an AHRS and on the coordinate frame to which the smartphone device (and the related measures) is referred. Subsection III-B presents in a concise but detailed way the deep architecture models analysed and tested during the work.

A. ORIENTATION ESTIMATION OVERVIEW

The orientation definition for a rigid body is generally made through a transformation matrix containing a parametrization of the Euler angles, unit quaternions, rotation vectors or rotation matrices [61]. Among them, the Euler angles allow for a more intuitive analysis in the 3D space and can be defined as follows:

- ϕ is known as *roll* angle and defines the x axis rotation;
- θ (*pitch* angle) refers to the y axis rotation;
- ψ (*yaw* or *heading* angle) represents the z axis rotation.

The correct integration of the raw IMU data or of the more cost-effective AHRS is at the basis of the orientation estimation process. The accelerometer measures the acceleration in m/s^2 applied to a device, including the force of gravity: velocity is determined if the linear acceleration component is integrated once and position if the integration is performed twice. The results can be of poor accuracy due to the extensive noise and accumulated drift from which it suffers. The rotation angles can be obtained by the integration of the angular velocities in rad/s

provided by the gyroscope; even if they are sensible to sudden and fast motion, these sensors generally experience major drift issues due to the errors accumulation over long time. For the aforementioned reasons, pose estimation is usually exploited through gyroscopes and accelerometers fusion to leverage their potential whilst attenuate their weaknesses. The Earth's magnetic field (μT) measures provided by the magnetometer can be joined to the previous ones to improve the heading determination; however, they suffer from the influence of metallic objects, which can heavily impact on the accuracy of the data collection. Moreover, the overall drift introduced by the sensors system causes errors accumulation: this means that the navigation information reliability and accuracy are guaranteed only within short times, with their measurements precision decreasing throughout long missions. For this reason, the integration of the measurements provided by the three sensors aims at reducing the errors accumulation caused by the single one; this is generally made through filtering techniques and fusion methods. Moreover, information provided by external devices can considerably improve the accuracy of the estimations, especially when low-cost sensors could facilitate the process and make it more practical.

In this context, the objective of the present work is to provide a supportive mean to improve the attitude estimations obtained by common AHRS: DOES is a low-cost DL architecture developed to recover orientation information from the view of a camera pointing the horizon at sea, which will be placed on the bow of a navigating vehicle in future experiments. The training has been performed on the ROPIS dataset, acquired using an application developed for the scope on an Android smartphone which simultaneously collects the frames and calculates the corresponding Ground Truth data using the AHRS sensors.

The IMU-AHRS measurements of the smartphones are generally expressed in a custom body reference frame. The Android developer website defines its frame relative to the device's screen when the device is held in its default orientation (see Fig. 2, [62]). In particular, the frame originates in the center of the device with the horizontal x axis pointing to the right, the vertical y axis pointing up and the z axis points toward the outside of the screen face, so that the coordinates behind the screen have negative Z values. The related attitude information is then referred to the same coordinates.

During the ROPIS dataset acquisition the smartphone has been kept in landscape mode, recording the horizon view. It has to be noticed that the coordinate frame does not change its definition, so in this setting the z axis points in the user direction, the y axis to his/her left and the x upwards.

B. DEEP LEARNING ARCHITECTURES

DOES model is composed of a pre-trained backbone CNN and two additional Fully Connected (FC) layers to output the roll and pitch estimates. Several, well established architectures have been tested as backbone for the final

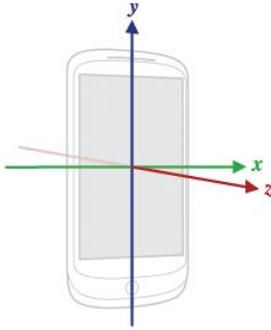


FIGURE 2. Device coordinate system used by the Android Sensor API [62].

network, as for example the VGG16-19 [63] and ResNet18-50-152 [64]; the resulting numerical comparison will be reported in Section VI, Table 3.

The VGG-16 and VGG-19 networks are based on the popular VGG architecture. They are composed of several convolutional layers followed by a Rectified Linear Unit (ReLU) activation function and interspersed by max pooling layers. Two FC layers are concatenated in order to produce the final features which are fed to a classification layer. These two networks differ only by the quantity and dimension of the convolutional layers employed, with a total number of parameters equal to 138M and 144M respectively. Despite being among the first developed deep architectures, with a huge amount of trainable parameters making them prone to overfitting, VGG models are still incredibly widespread, thanks to their ease of use for fine-tuning purposes on different tasks [65], [66].

ResNet is a family of deep models based on the residual architecture. Differently from the VGG, the ResNet is made of a series of residual blocks in which the feature maps calculated by the convolutional layers are added to the input, so that each residual block calculates an update (hence residual) of the input feature maps. This approach makes the network resilient to the vanish gradient problem [67], improving convergence speed and the final accuracy result. Moreover, all the ResNet models avoid the use of the FC layers after the convolutional blocks, reducing the total number of trainable parameters and thus lessening the overfitting effect on training data. Authors of ResNet developed three versions with different number of layers (18, 50, 152) and with different number of visual features before the classification step (512 for the former, 2048 for the others). The number of free parameters for the 18, 50 and 152 layers models are 11M, 23M and 60M respectively.

In the experiments presented in this work, all the networks have been fine-tuned on the proposed ROPIS dataset starting from the ImageNet [68] pre-trained weights. The ResNet18 has been chosen among the others as the default DOES backbone since it produced the best accuracy while keeping at the same time a fast inference speed. Fig. 3 reports the DOES network with the default ResNet18 backbone.

Two additional FC layers have been added as additional branches on top of the highest set of visual features in the backbone network to separately estimate the roll and pitch angles; for example, in the case of the ResNet models, this correspond to the global average pooling layer. Some different estimation procedures have been experimented, as the one described in [69]: it proposes to map the float angle value to a set of fixed bins, which then undergo a standard classification procedure with a final mapping back to the float value. However, in this work it has been experimentally found that this approach adds a layer of complexity without increasing the overall performances; this led to the decision to add a FC layer for each angle, which is able to accomplish the regression task with a good accuracy. Both the backbone network and the additional FC layers are jointly trained by back-propagation with the use of a standard Mean Square Error Loss (squared L2 norm). Two separated losses are calculated for each of the two angles as reported in (1) for roll (L_{roll}) and (2) for pitch (L_{pitch}), where y and \hat{y} are the GT and predicted values respectively. The final loss L_{final} is then obtained as a simple addition of the aforementioned quantities, as shown in (3). The GT roll and pitch values have undergone a prior normalization process, which subtracts to each of them the mean and divides by the variance, both calculated over the entire dataset.

$$L_{roll}(y_{roll}, \hat{y}_{roll}) = \frac{1}{n} \sum_{i=1}^n (y_{roll} - \hat{y}_{roll}^i)^2 \tag{1}$$

$$L_{pitch}(y_{pitch}, \hat{y}_{pitch}) = \frac{1}{n} \sum_{i=1}^n (y_{pitch} - \hat{y}_{pitch}^i)^2 \tag{2}$$

$$L_{final} = L_{roll}(y_{roll}, \hat{y}_{roll}) + L_{pitch}(y_{pitch}, \hat{y}_{pitch}) \tag{3}$$

IV. ROPIS DATA ACQUISITION PROCESS

The lack of datasets designed for DL-based orientation estimation at sea lead to the necessity of searching for methods to acquire a set of data for the scope. In the following section, the development of the Android application and the obtained ROPIS dataset will be described in detail.

A. DEVICE INTERNAL SENSORS AND CHARACTERISTICS

In order to train the model, the dataset needs to contain a large amount of images showing the horizon and the corresponding GT data in terms of roll and pitch angles. The latter needs to be given with the best possible accuracy, as the learning process results will depend on it, which is strictly related to the instrumentation employed for the acquisition. With the aim of producing a low-cost and flexible solution, in this work the authors avoided the use of costly, high-end IMU devices and developed the FrameWOAndroid application to acquire the dataset through a common smartphone. The presented ROPIS dataset in its first release has been totally collected through a OnePlus Nord smartphone, equipped with the most common sensors (Table 1) and characterized by an average price.

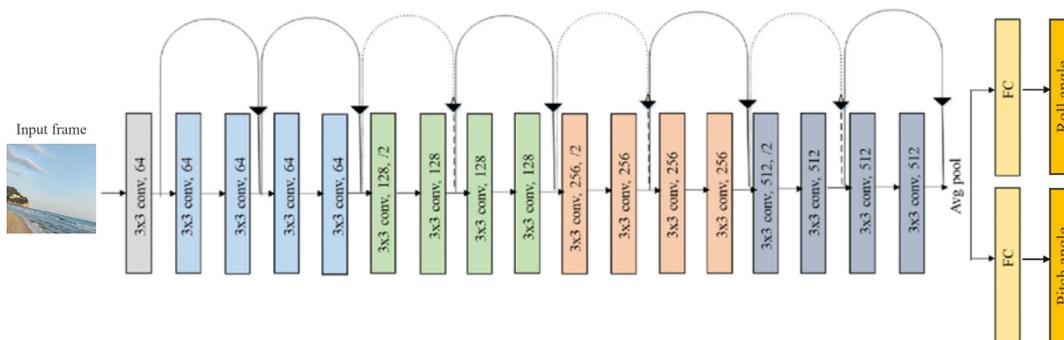


FIGURE 3. DOES architecture with default ResNet18 backbone network.

TABLE 1. OnePlus Nord smartphone general specifics [72].

General	Main Sensors	Rear Camera - Main
OS: OxygenOS Android™ 10	IMU: Bosch BMI260	Megapixels: 48
CPU: Qualcomm® Snapdragon™ 765G	Magn: MEMSIC MMC5603	Pixel Size: 0.8 μm/48M; 1.6 μm (4 in 1)/12M
GPU: Adreno 620	Camera: Sony IMX586	Lens Quantity: 6P
RAM: 8GB/12GB LPDDR4X	Proximity sensor	Aperture: f/1.75
Storage: 256GB UFS2.1	Ambient light sensor	OIS, EIS: Yes

The OnePlus Nord mounts a BMI260 IMU, which contains a 16-bit tri-axial gyroscope (G) and accelerometer (A) providing fast, precise inertial sensing in smartphones and Human-Machine Interface (HMI) applications (i.e., advanced gesture, activity and context recognition, etc.). The IMU is characterized by a noise density of $160\mu g/\sqrt{Hz}$ (A) and $0.008 dps/\sqrt{Hz}$ (G), a Zero-g/Zero-rate offset of $\pm 20 mg$ (A) and $\pm 0.5 dps$ (G) and an output data rate up to $1.6 kHz$ (A) and $6.4 kHz$ (G). Moreover, it mounts the industry’s first self-calibrating gyroscope with motionless Component Re-Trimming (CRT) functionality, which compensates MEMS typical soldering drifts, ensuring post-soldering sensitivity errors down to $\pm 0.4\%$ [70].

The MMC5603 is a monolithic complete 3-axis Anisotropic Magnetoresistance Effect (AMR) magnetic sensor. It has an on-chip automatic degaussing with built-in SET/RESET function which eliminates the thermal variation-induced offset error and clears the residual magnetization deriving from strong external fields. Its true frequency response is up to $1kHz$ and can measure magnetic fields in a range of $\pm 30Gauss$ (G) with $2mG$ total Root Mean Square (RMS) noise level, enabling heading accuracy of $\pm 1deg$ in electronic compass applications [71].

The Sony IMX586 stacked CMOS image sensor is mounted as the main camera of the OnePlus Nord, and features 48 effective megapixels with an ultra-compact pixel size of $0.8\mu m$. The sensor uses the Quad Bayer color filter array, where adjacent 2×2 pixels come in the same color, making high-sensitivity shooting possible. During low light shooting, the signals from the four adjacent pixels are added, raising the sensitivity to a level equivalent to that of $1.6\mu m$ pixels (12 megapixels), resulting in bright, low noise images [73].

B. FrameWO APPLICATION DEVELOPMENT

The FrameWO app has been developed in a free Open Source environment, the B4X suite [74], which supports the majority of PC, smartphones and embedding operating systems (e.g., Android, iOS, Windows, MacOS, Linux, Arduino, RaspberryPI) and uses a modern version of Visual Basic as programming language. The Android version (B4A) allows to wrap existing Java code as an external library and then to reference it from the B4A IDE, obtaining in release mode performances similar to those of Java. The size of a simple app is generally around 100 KB.

As previously mentioned, the necessary prerequisite for the dataset to meet the scope of this study is to associate to each frame the corresponding GT; however, the images size is much more larger than that of the IMU data, thus introducing a delay in their storage which affected their simultaneity. For this reason, the app captures the frames in YUV format (allowing for a better compression of the image) and converts them in JPEG only at the end of the process; this also avoids to run out of memory during the acquisition. A detailed overview on the YUV model can be found in [75]. Furthermore, several tests have been performed to determine an acquisition frequency value suitable for both the high-rate IMU data and the low-rate camera frames: the application offers in fact the possibility to set the camera acquisition frequency in msec to choose the best option for the needs.

As regards the GT, the API of Android [62] has been used to work on the raw measures read by the sensors and to obtain the Euler angles of interest. The *getRotationMatrix* function allows for a coordinate systems transformation (from the device to the world one in this case) and takes as input the gravity and geomagnetic field in vector form to compute the inclination matrix *I* and the rotation matrix *R*.

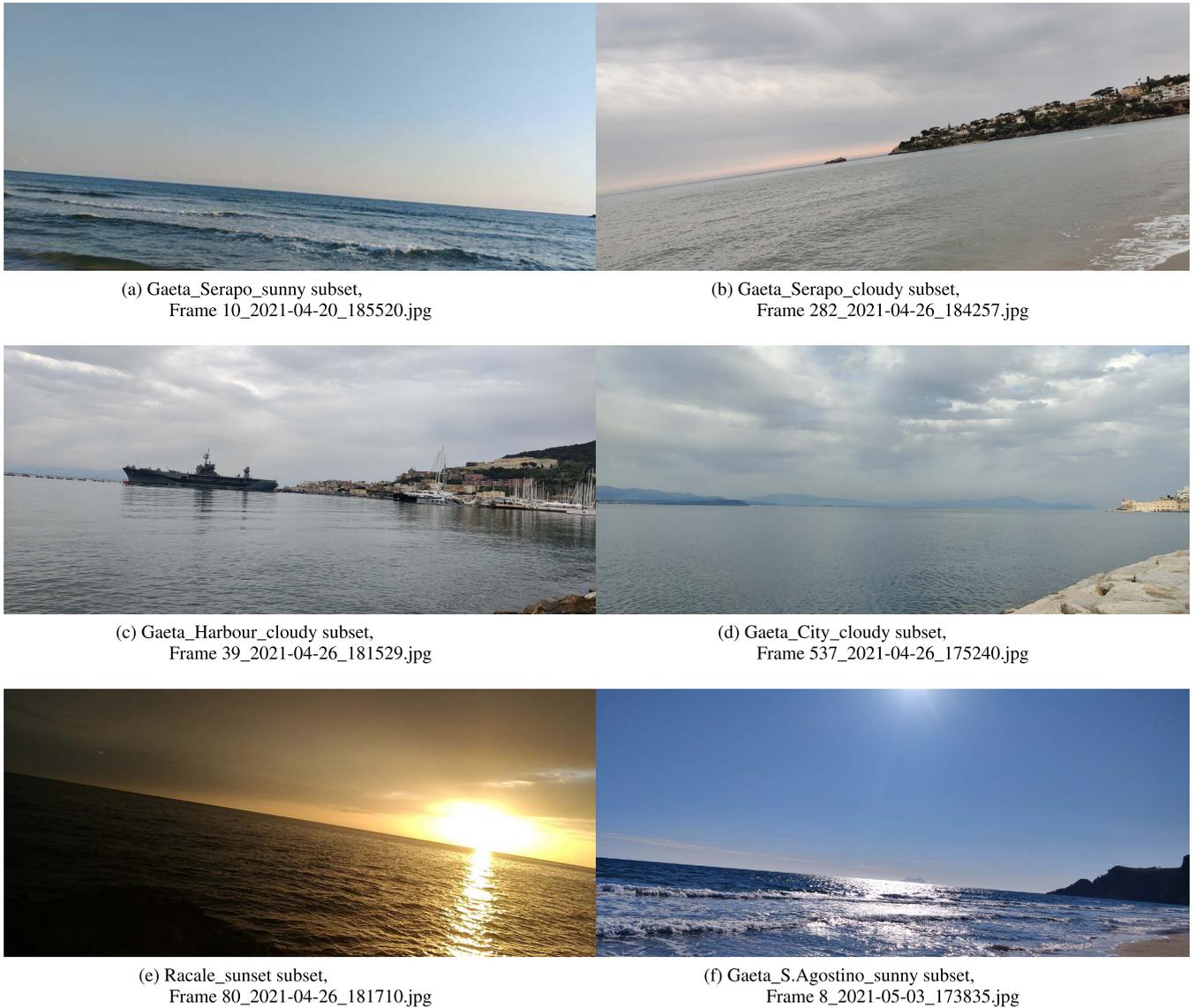


FIGURE 4. ROPIS dataset samples. Fig. 4a to 4e belong to the training set, Fig. 4f to the test set.

By definition, I is the rotation around the X axis which converts the geomagnetic vector into the gravity coordinate space, whereas R defines the identity matrix of the device aligned with the same world coordinate system: in this setting, the device faces the sky with the X axis pointing the East and the Y axis the North Pole (see (4), where g is the magnitude of gravity and m is the magnitude of the geomagnetic field).

$$\begin{bmatrix} 0 & 0 & g \\ 0 & m & 0 \end{bmatrix} = I * R * \text{geomagneticfield} \quad (4)$$

In order to isolate the gravity vector, a discrete-time low-pass filter with a smoothing factor $\alpha = 0.2$ has been applied to the accelerometer measurements. The Euler angles are recovered through the *getOrientation* function, which calculates them from the elements of the rotation matrix R [62], [76].

The measurements are updated at the fastest rate provided by the Android API, which is in the order of few milliseconds. The time sampling has been set equal to 100msec, that means that 10 times in a second the device simultaneously registers the orientation and the corresponding image. As a final result, data are saved in a directory named with the date and time of the specific acquisition, which is further renamed to specify the scenario characteristics of the moment. This directory contains all the frames, saved as n_YYYY-MM-DD_HHMMSS.jpg, and a data.txt file which lists the frame name, its index n , and the related GT.

C. DATASET STRUCTURE

The ROPIS dataset in its first release has been mainly acquired in Italy, in the cities of Gaeta (Lazio) and Racale (Puglia). It consists of 22173 sRGB TrueColor JPEG images, with resolution set to 2592×1168 , for a total dimension

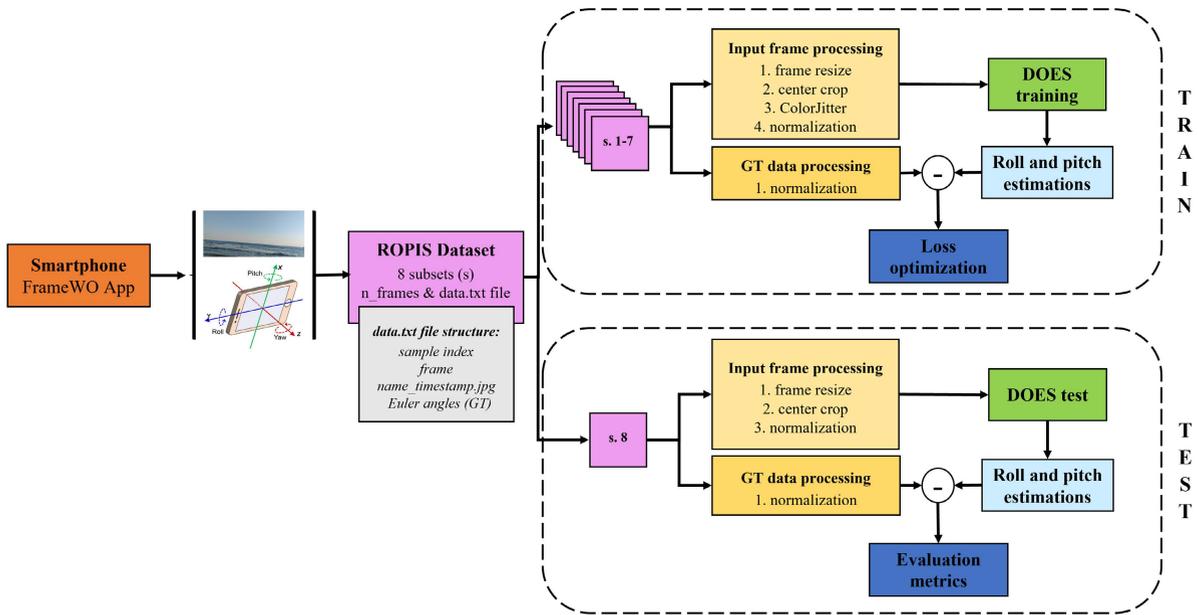


FIGURE 5. Detailed DOES workflow, showing the data acquisition and the train/test phases.

of 42.3 GB. Six different subsets have been acquired in as many locations, each presenting different characteristics in terms of scenarios and meteo-marine conditions; five of them have been chosen for the training set, from which a total of 100 frames has been separated for the validation set, and the last acquisition has been used as test set. The use of a dedicated test set with images coming from a separate location allows to verify the ability of DOES to generalize to new, different scenes with respect to the training and validation set. More in the specific, in each place eight different acquisitions have been made trying to simulate the behaviour of a ship in navigation in both static and dynamic conditions: this aims at emulating the induced oscillations which resemble the true motion of the ship. To improve the generalization ability of the model, the data have been acquired at different day times and with sunny and cloudy sky; Fig. 4 shows different samples of the ROPIS dataset. Some aspects of these data need to be highlighted:

- The point of view of the ROPIS images presents some differences with respect to the acquisitions taken on board the ship, since it adds parts of the land in the image foreground, such as sand, rocks, etc. However, this does not affect the learning procedure as the DL networks are able to recognize useful and useless image features, discarding the latter.
- A frame representing the real view from a navigating vehicle should depict some elements in the scene, such as the bow structures and some part of the bridge floor from a ship, or some of the USV sections. Although these specific features do not appear in ROPIS, DOES demonstrated its robustness to similar images cluttering present in the frames. Further experiments will be made to precisely assess their impact on the learning process.

- The data acquisition has been made with the camera at a roughly fixed height of 1.5m with slight oscillations around this value: this considers, among the different vehicle movements, also the linear vertical -up/down-motion along the z axis (*heave*), corresponding to the smartphone x axis. It should be remarked that the pitch estimation is strictly related to the horizon height and thus to the camera axis and view; for this reason, the horizon line should be obviously always visible in the frame.

Fig. 5 shows the workflow of DOES in its three main phases: the data acquisition, the training with its specific data augmentation process and the test which finally allows to calculate the evaluation metrics.

The ROPIS dataset is intended to be further enhanced. The use of other low-cost cameras (to take into account the differences in the camera parameters and lens distortion) and the setting of a range of different camera height values aim at considering their impact on the training phase. Moreover, the acquisitions will be made in different scenarios, which will include adverse meteo-marine conditions and locations as ships bridge and USV platforms. The heterogeneity of the data fed to the network will enhance the model capability to generalize over more complex data and realistic settings, making it invariant to these parameters.

V. EXPERIMENTAL SETUP

In this section some details on the training process will be given, together with a brief overview of the evaluation metrics used to appraise the performance of DOES. Finally, the problem related to the comparison of DOES with other methods will be discussed.

A. TRAINING DETAILS

DOES has been developed in Python programming language using the Pytorch framework; the code is publicly available.¹ DOES has been trained using a standard fine-tuning procedure: the backbone convolutional kernels were pre-trained on ImageNet while the additional FC layers have been initialized with random values drawn upon Pytorch default uniform distribution. Both convolutional and FC layers have been trained using the Adam optimizer [77] and a fixed learning rate set to 0.001. DOES has been trained on the ROPIS training set for a total of 10 epochs: it has in fact been noticed that a larger number of epochs led to an increase of the overfitting without any improvement of the accuracy.

The images have been squared to a preliminary 2592×2592 resolution by the application of a zero-padding; this operation adds black bands to the smallest dimension to obtain a squared input whilst preventing the loss of information. The images have then been resized to a final resolution of 224×224 ; a zero mean-unit variance normalization has been applied to both the images and the GT sets, with the corresponding mean and variance calculated over the specific training data.

The data augmentation process consisted of random changes in the colours of the images, using the *ColorJitter* transformation function of Pytorch which allows to set different values of brightness, contrast, saturation and hue: this resulted in an increase of the training dataset which further enhanced the generalization abilities of DOES. No random cropping nor image flipping have been applied during this process: in fact, the former would have caused the neglecting of the relative sea height information given by the images whereas the latter could have changed the correct roll angle perception of the network. The data augmentation procedure has naturally been deactivated during the testing phase, whereas the zero-padding and resize processes have been applied also to the test images; furthermore, the predicted roll and pitch values have been de-normalized before calculating the evaluation metrics presented in the following paragraph V-B. The selected data augmentation values (brightness and hue equal to 0.5, contrast and saturation equal to 5), as well as all the other training hyperparameters, have been tuned on the validation set.

B. EVALUATION METRICS

DOES has been evaluated on the basis of the regression metrics implemented by the Scikit library in the *sklearn.metrics* module, which contains the most common utility functions to measure the regression performance.

The Mean Absolute Error (MAE) computes a risk metric corresponding to the expected value of the absolute error (5); it is the average absolute difference between the predicted and the true value, expressed in the same scale as the data being measured. Each error contributes to MAE in proportion to its

absolute value.

$$MAE(y, \hat{y}) = \frac{1}{n} \sum_{i=0}^{n-1} |y_i - \hat{y}_i| \quad (5)$$

The Root Mean Square Error (RMSE) represents the square root of the second sample moment of the differences between predicted values and the observed values (or the quadratic mean of these differences, also called residuals). It is a measure of accuracy and it is sensitive to outliers (6). In fact, since the errors are squared before they are averaged, the RMSE gives a relatively high weight to large errors, making it more useful when large errors are particularly undesirable. RMSE does not necessarily increase with the variance of the errors, growing instead with the variance of the frequency distribution of error magnitudes.

$$RMSE(y, \hat{y}) = \frac{1}{n} \sqrt{\sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2} \quad (6)$$

The Standard Deviation (STD) is a measure of the amount of dispersion (or variation) of the samples. A low standard deviation indicates that the values tend to be close to the mean μ (also called the expected value) of the set, whereas a high standard deviation indicates that the values are spread out over a wider range (7).

$$\sigma(\hat{y}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \mu)^2} \quad (7)$$

Finally, the Median Absolute Error (MedAE) is calculated by taking the median of all the absolute differences between the GT and the prediction (8). It is a non-negative floating point with best value of 0.0, robust to outliers since the median is not affected by values at the tails.

$$MedAE(y, \hat{y}) = \text{median}(|y_i - \hat{y}_i|, \dots, |y_n - \hat{y}_n|) \quad (8)$$

C. METHODOLOGY COMPARISON

The comparison between DOES and other state of the art methods turned out to be a non trivial task for several reasons; among the others, the Deep Learning based solutions currently developed for the estimation of roll and pitch are either released without source code (as for example in [35]) or employed for very different tasks (e.g., head pose estimation [69]), thus making the comparison not properly correct or practically impossible. Generally speaking, traditional Horizon Line Detection (HLD) algorithms can be used as a proxy for this kind of estimations; the roll and pitch angles can in fact be correlated to the slope and position of the horizon line. However, as previously mentioned, this would require the correct knowledge of the intrinsic and extrinsic camera parameters and of the transformation matrix between the camera and the smartphone reference systems. To address this problem, a Linear Least Squares method has been applied to calibrate the HLD algorithms on the basis

¹<https://github.com/fabidicia/does>

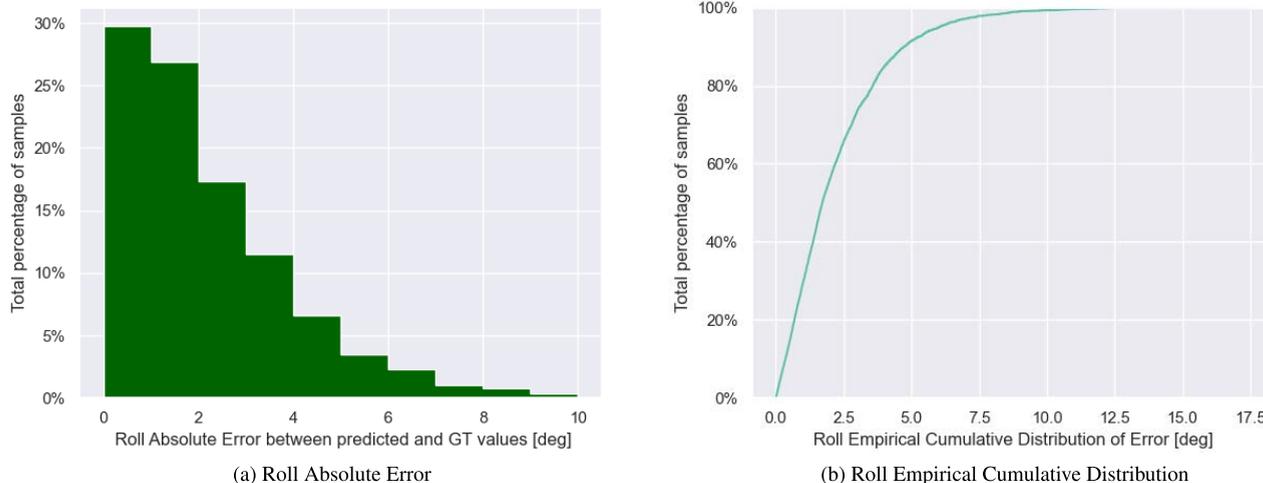


FIGURE 6. Graphical distribution of the errors for the estimation of the roll angle.

of the minimization of the squared error calculated between their output predictions and the GT values.

Two of the most renowned HLD algorithms by the scientific community have been selected to perform this comparison and are briefly described in the following lines.

The **Otsu** method [78] is a popular technique used to threshold the image between sky and non-sky regions. It is a reasonable fast and simple algorithm which performs fairly well on heterogeneous sets of data. The threshold value T is automatically computed by the algorithm through the assumption that the grayscale histogram of the image pixels intensities is bi-modal; the threshold is set so that the distance between the two histogram peaks is maximized.

Ettinger et al. [79] is a computer vision-based HLD algorithm that performs exhaustive search in the 2D line parameters space over the whole image looking at the best values which separate sky from terrain. However, being a slow algorithm on high resolution images, a modified version has been implemented that uses a two-stage objective: the *global* one searches for a narrow range of combinations of the pitch and roll horizon line angles corresponding to a half-plane that likely subdivides the sky from the rest of the image. The *local* one aims at searching exhaustively through these combinations to find the half-plane that maximizes the difference (in average intensity) of the two half-planes in their immediate vicinity. This method assumes that the sky pixels have higher intensity values than the ground pixels (higher mean), and that the sky has higher consistency of representation (lower variance).

VI. RESULTS AND DISCUSSION

This section contains an assessment of the results provided by DOES. Table 2 shows DOES performances with respect to the selected horizon line detection algorithms. DOES is able to achieve sensible better results both on roll and pitch

TABLE 2. DOES performances compared to those of the two HLD methods.

	DOES		Otsu [78]		Ettinger [79]	
	<i>roll</i>	<i>pitch</i>	<i>roll</i>	<i>pitch</i>	<i>roll</i>	<i>pitch</i>
MAE [deg]	1.65	1.84	4.48	3.76	4.04	3.77
RMSE [deg]	2.27	2.45	5.44	4.75	5.01	4.78
STD [deg]	1.55	1.61	3.09	2.90	2.97	2.93
MedAE [deg]	1.14	1.41	4.04	3.19	3.44	3.15

angles, with a Mean Absolute Error close to 1.5° , as opposed to the other methods which exhibit worse performance on all the indicators.

The MAE and the RMSE can be used together to diagnose the variation in the errors in a set of predictions. The RMSE is generally higher than the MAE, and the greater is the difference between them, the greater will be the variance in the individual errors of the samples; moreover, if the RMSE is close to the MAE, then all the errors are of the same magnitude. In the case of the current comparison, the small gap between RMSE and MAE demonstrates the ability of DOES to produce fewer outliers than Otsu and Ettinger. In addition, the STD values of the three methods show that the results obtained by DOES are significantly more clustered than the others, meaning that they are closer to the mean value and as such can be considered more reliable. The good performances of DOES are further confirmed by the MedAE value, which is sensibly lower than the counterparts. These findings can be summarized in Fig. 6, which shows the MAE behaviour analysing the outputs percentage belonging to different MAE intervals (Fig. 6a) together with the empirical cumulative distribution (Fig. 6b) for the roll angle. The same evaluation can be made for the pitch angle (Fig. 7), which exhibits similar performances to the roll angle. Another important consideration related to this comparison regards the inference time of DOES; the average

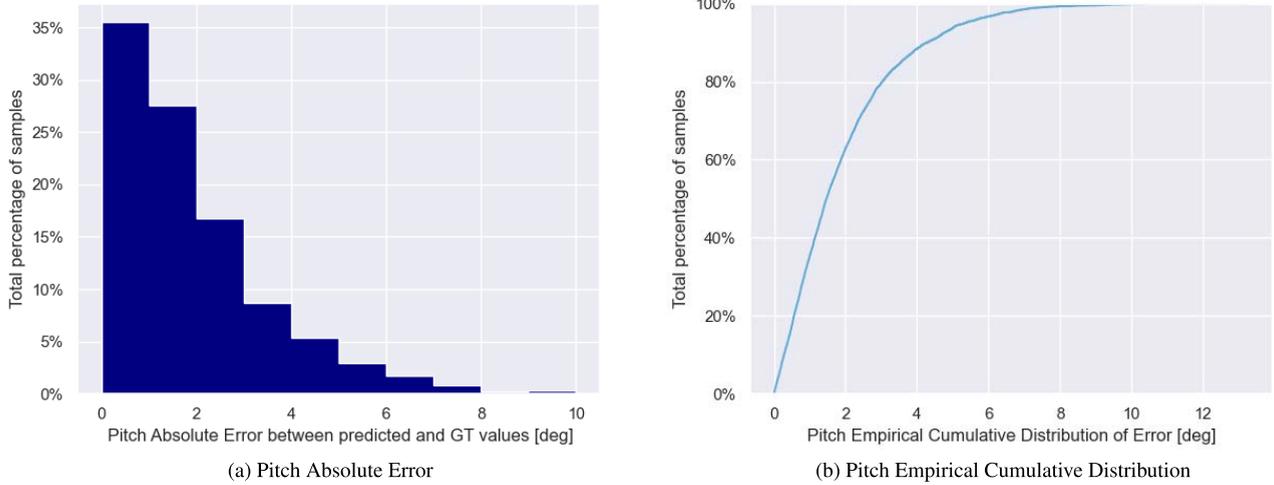


FIGURE 7. Graphical distribution of the errors for the estimation of the pitch angle.

TABLE 3. Comparative results on different DOES backbones. TP indicates the number of trainable parameters.

	ResNet18 TP = 11M		ResNet50 TP = 23M		ResNet152 TP = 58M		VGG19 TP = 139M		VGG19bn TP = 139M		DenseNet161 TP = 26M	
	roll	pitch	roll	pitch	roll	pitch	roll	pitch	roll	pitch	roll	pitch
MAE [deg]	1.65	1.84	1.77	1.88	1.82	1.92	4.67	4.11	1.91	1.99	1.63	1.87
RMSE [deg]	2.27	2.45	2.40	2.51	2.44	2.54	5.60	5.18	2.57	2.61	2.23	2.48
STD [deg]	1.55	1.61	1.63	1.66	1.62	1.66	3.09	3.14	1.72	1.69	1.55	1.63
MedAE [deg]	1.14	1.41	1.28	1.46	1.36	1.52	4.26	3.43	1.47	1.56	1.14	1.44

estimation time on a single image is 100-150msec with any of the tested backbones, whereas Otsu and Ettinger inference time is comprised between 100 and 11000 msec, making them unsuitable for real-time applications on high-resolution images.

Table 3 shows a detailed comparison between DOES with its default proposed network and some alternative backbones: DOES is able to produce good performances with all the residual networks, whereas both VGG-19 and VGG-19bn struggle to produce reasonable results. More in detail, the MAE and RMSE results of ResNet18 are slightly better than the 50- and 152-layers versions, with the powerful DenseNet161 model able to produce a similar accuracy only on the roll angle. The performing results obtained by the ResNet18, together with the fastest training and inference speed (due to the smaller number of trainable parameters TP with respect to the other architectures), make ResNet18 the first choice for the deployment of DOES as long as new models specifically developed for the scope will be released. Future work will focus on the use of lighter architectures developed for the specific use on low-resources embedded hardware (e.g., MobileNet, [80]); this will lay the foundation for the deployment of the proposed model on embedded devices (e.g., Nvidia Jetson, [81]) in real-time scenarios, in accordance with the aim of making DOES a supportive smart technology to improve the attitude estimations provided by low-cost sensors.

Furthermore, the ROPIS dataset has been used for an additional test in which a 1.33x zoom has been applied to the frames to simulate different camera parameters. In some cases, this corresponded to a crop in the image which removed the horizon line, thus making DOES unable to correctly estimate the angles. This reflects in a slight decrease of the performances: the roll MAE is equal to 2.10°, with a RMSE of 2.81°, whereas the pitch angle exhibits a 2.02° MAE and a 2.90° RMSE.

Finally, a separated test (with no prior training or specific tuning) has been made on a set of 191 images presenting three main variations with respect to the ROPIS train and test data:

- The device: a smartphone Huawei P9 [82] has been used, with the FrameWO App, to collect the data. The mounted dual-lens Leica camera has different characteristics with respect to the OnePlus Nord Sony camera: the P9 Leica 12 MP has in fact an aperture size of $f/2.2$, a focal length of 27mm (wide), a sensor size of $1/2.9''$ and a pixel size of $1.25\mu\text{m}$.
- The location: the acquisition has been made in a different area of the Racale city (LE).
- The environment setting: the data have been collected rightly after the sunset, in a low-light condition which highly reduced the contrast in the frame, resulting in a very challenging scenario.

Despite these substantial changes in the sensor and in the overall acquisition, DOES obtained remarkable results,

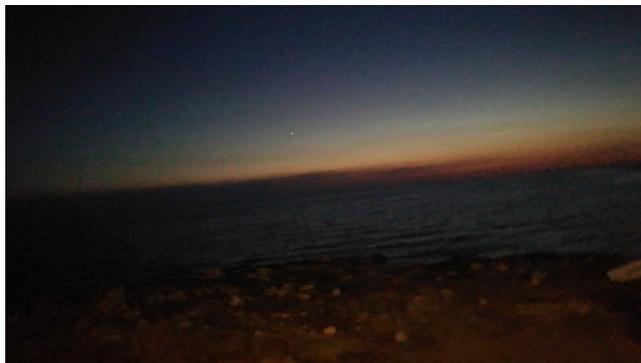


FIGURE 8. A frame from the low-light condition separated set.

performing a 2.17° MAE and a 2.70° RMSE for the roll angle and a 2.22° MAE and a 2.71° RMSE for the pitch angle. This demonstrates that DOES can successfully generalize over various conditions and camera parameters, confirming its potential for more challenging settings and further employment as inertial systems support and visual-based odometry tasks.

It is worth mentioning that the accuracy of the results is proportioned to the precision of the GT data and thus of the systems employed to acquire it. In this case, the overall accuracy is strictly connected to the use of a smartphone AHRS which, although being limited to the low-cost sensors mounted on it, is still able to provide reliable and accurate measurements. The use of high-end and more expensive devices would in fact ensure a higher grade of GT accuracy with consequent improvements in the DOES performances.

VII. CONCLUSION

This paper presents a novel Deep Learning-based approach to the attitude estimation problem, which has been developed and intensively tested on a new dataset (the ROPIS dataset) specifically built for the scope and released in the context of this work. Deep Orientation (of roll and pitch) Estimation at Sea (DOES) is able to predict the attitude of the device in terms of roll and pitch angles by analysing the frames recorded by the camera pointing towards the sea horizon. DOES has been tested using several known architectures (e.g., ResNet152, ResNet18, VGG19) and with different configurations and hyper-parameters, obtaining excellent results. Unlike other visual-based methods, DOES is able to produce the output without the explicit knowledge of the camera intrinsic and extrinsic parameters or the distortions introduced by the camera lens. There is in fact no necessity to make any assumption on the use of specific models to parametrize the camera, since the model training only depends on the dataset given as input; the latter generally provides different sampling characteristics, thus making the network able to learn and then estimate the attitude regardless of the camera specifics.

The ROPIS dataset has been created for this particular task and is here presented in its first release; the lack of public

datasets suitable for DL applications made it necessary to search for a valid alternative for the experiments conduction. For this reason, the FrameWO Android application has been developed using the Open Source B4A platform and will be made publicly available online. This app allows to simultaneously acquire the frames to be fed to the model as input, and the attitude estimations measured through the internal sensors of the smartphone, which will be used as Ground Truth in the training/testing phases.

ROPIS dataset is intended to be further improved by the introduction of more subsets of data collected in different scenarios (i.e., during the dusk/dawn, rainy days, etc) and environments (e.g., different cities coastlines, onboard of a vessels), using different acquisition devices. This will improve the DOES ability to generalize over heterogeneous data, making it even more invariant to the camera configurations, the acquisition condition and cluttering factors, thus providing better results in any kind of situation in which the vehicle will be navigating. In this regard, the authors wish to encourage the users to download and test the FrameWO application with the aim of enhancing the ROPIS and its usage among the scientific community, to give a concrete contribution to this task.

The objective of this project is to develop a supportive technology to be integrated to the existing low-cost methodologies employed for the attitude estimation task. In fact, it has to be noticed that this approach has been specifically designed using affordable devices and applications and, as such, its results are not intended (at least in its preliminary version) to reach the accuracy provided by high-precision modern sensors. Further experiments will be made to test other light-weight DL architectures, which could be deployed on low-resources embedded hardware with the aim of providing better accuracy results in real-time applications on autonomous vehicles. These enhancements will make DOES a robust system to be integrated in visual and visual-inertial odometry methodologies.

ACKNOWLEDGMENT

The authors would like to give a special thanks to Mr. Alberto Greco, which has been of fundamental importance in the development stage of the app employed to acquire the ROPIS dataset.

REFERENCES

- [1] M. B. Alataise and G. P. Hancke, "Pose estimation of a mobile robot based on fusion of IMU data and vision data using an extended Kalman filter," *Sensors*, vol. 17, no. 10, p. 2164, 2017.
- [2] S. Del Pizzo, S. Gaglione, A. Angrisano, G. Salvi, and S. Troisi, "Reliable vessel attitude estimation by wide angle camera," *Measurement*, vol. 127, pp. 314–324, Oct. 2018.
- [3] U. Ganbold and T. Akashi, "The real-time reliable detection of the horizon line on high-resolution maritime images for unmanned surface-vehicle," in *Proc. Int. Conf. Cyberworlds (CW)*, Sep. 2020, pp. 204–210.
- [4] G. Huang, "Visual-inertial navigation: A concise review," in *Proc. Int. Conf. Robot. Automat. (ICRA)*, May 2019, pp. 9572–9582.
- [5] L. Han, Y. Lin, G. Du, and S. Lian, "DeepVIO: Self-supervised deep learning of monocular visual inertial odometry using 3D geometric constraints," 2019, *arXiv:1906.11435*.

- [6] S. Poddar, R. Kottath, and V. Karar, "Evolution of visual odometry techniques," 2018, *arXiv:1804.11142*.
- [7] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza, "On-manifold preintegration for real-time visual-inertial odometry," *IEEE Trans. Robot.*, vol. 33, no. 1, pp. 1–21, Feb. 2017.
- [8] M. Kok, J. D. Hol, and T. B. Schön, "Using inertial sensors for position and orientation estimation," *Found. Trends Signal Process.*, vol. 11, nos. 1–2, pp. 1–153, Nov. 2017, doi: [10.1561/20000000094](https://doi.org/10.1561/20000000094).
- [9] A. Avci, S. Bosch, M. Marin-Perianu, R. Marin-Perianu, and P. Havinga, "Activity recognition using inertial sensing for healthcare, wellbeing and sports applications: A survey," in *Proc. 23th Int. Conf. Archit. Comput. Syst.*, 2010, pp. 1–10.
- [10] H. J. Luinge and P. H. Veltink, "Measuring orientation of human body segments using miniature gyroscopes and accelerometers," *Med. Biol. Eng. Comput.*, vol. 43, no. 2, pp. 273–282, Mar. 2005.
- [11] H. Fourati, N. Manamanni, L. Afilal, and Y. Handrich, "A nonlinear filtering approach for the attitude and dynamic body acceleration estimation based on inertial and magnetic sensors: Bio-logging application," *IEEE Sensors J.*, vol. 11, no. 1, pp. 233–244, Jan. 2011.
- [12] S. Adler, S. Schmitt, K. Wolter, and M. Kyas, "A survey of experimental evaluation in indoor localization research," in *Proc. Int. Conf. Indoor Positioning Indoor Navigat. (IPIN)*, Oct. 2015, pp. 1–10.
- [13] H. G. de Marina, F. J. Pereda, J. M. Giron-Sierra, and F. Espinosa, "UAV attitude estimation using unscented Kalman filter and TRIAD," *IEEE Trans. Ind. Electron.*, vol. 59, no. 11, pp. 4465–4474, Nov. 2012.
- [14] E. Vertzberger and I. Klein, "Attitude adaptive estimation with smartphone classification for pedestrian navigation," *IEEE Sensors J.*, vol. 21, no. 7, pp. 9341–9348, Apr. 2021.
- [15] V. Renaudin and C. Combettes, "Magnetic, acceleration fields and gyroscope quaternion (MAGYQ)-based attitude estimation with smartphone sensors for indoor pedestrian navigation," *Sensors*, vol. 14, no. 12, pp. 22864–22890, 2014.
- [16] R. Harle, "A survey of indoor inertial positioning systems for pedestrians," *IEEE Commun. Surveys Tuts.*, vol. 15, no. 3, pp. 1281–1293, 3rd Quart., 2013.
- [17] N. H. Q. Phuong, H.-J. Kang, Y.-S. Suh, and Y.-S. Ro, "A DCM based orientation estimation algorithm with an inertial measurement unit and a magnetic compass," *J. Universal Comput. Sci.*, vol. 15, no. 4, pp. 859–876, 2009.
- [18] A. Kim and M. F. Golnaraghi, "A quaternion-based orientation estimation algorithm using an inertial measurement unit," in *Proc. Position Location Navigat. Symp. (PLANS)*, 2004, pp. 268–272.
- [19] A.-J. BaerVELdt and R. Klang, "A low-cost and low-weight attitude estimation system for an autonomous helicopter," in *Proc. IEEE Int. Conf. Intell. Eng. Syst.*, Sep. 1997, pp. 391–395.
- [20] D. Gebre-Egziabher, G. H. Elkaim, J. D. Powell, and B. W. Parkinson, "A gyro-free quaternion-based attitude determination system suitable for implementation using low cost sensors," in *Proc. IEEE Position Location Navigat. Symp.*, Mar. 2000, pp. 185–192.
- [21] R. G. Valenti, I. Dryanovski, and J. Xiao, "Keeping a good attitude: A quaternion-based orientation filter for IMUs and MARGs," *Sensors*, vol. 15, no. 8, pp. 19302–19330, 2015.
- [22] B. Allotta, A. Caiti, R. Costanzi, F. Fanelli, D. Fenucci, E. Meli, and A. Ridolfi, "A new AUV navigation system exploiting unscented Kalman filter," *Ocean Eng.*, vol. 113, pp. 121–132, Feb. 2016.
- [23] W. Li and J. Wang, "Effective adaptive Kalman filter for MEMS-IMU/magnetometers integrated attitude and heading reference systems," *J. Navigat.*, vol. 66, no. 1, pp. 99–113, 2013.
- [24] F. Di Ciaccio, S. Gaglione, and S. Troisi, "A preliminary study on attitude measurement systems based on low cost sensors," in *Proc. Int. Workshop R3 Geomatics, Res., Results Rev.* Cham, Switzerland: Springer, 2019, pp. 103–115.
- [25] T. Michel, P. Geneves, H. Fourati, and N. Layaïda, "On attitude estimation with smartphones," in *Proc. IEEE Int. Conf. Pervasive Comput. Commun. (PerCom)*, Mar. 2017, pp. 267–275.
- [26] R. V. Vitali, R. S. McGinnis, and N. C. Perkins, "Robust error-state Kalman filter for estimating IMU orientation," *IEEE Sensors J.*, vol. 21, no. 3, pp. 3561–3569, Feb. 2021.
- [27] K. Wen, K. Yu, Y. Li, S. Zhang, and W. Zhang, "A new quaternion Kalman filter based foot-mounted IMU and UWB tightly-coupled method for indoor pedestrian navigation," *IEEE Trans. Veh. Technol.*, vol. 69, no. 4, pp. 4340–4352, Apr. 2020.
- [28] D. A. Aligia, B. A. Roccia, C. H. De Angelo, G. A. Magallán, and G. N. González, "An orientation estimation strategy for low cost IMU using a nonlinear Luenberger observer," *Measurement*, vol. 173, Mar. 2021, Art. no. 108664.
- [29] J. Schnee, J. Stegmaier, T. Lipowsky, and P. Li, "Auto-correction of 3D-orientation of IMUs on electric bicycles," *Sensors*, vol. 20, no. 3, p. 589, Jan. 2020.
- [30] B. Wang, Y. Su, and L. Wan, "A sea-sky line detection method for unmanned surface vehicles based on gradient saliency," *Sensors*, vol. 16, no. 4, p. 543, Apr. 2016.
- [31] C. Y. Jeong, H. S. Yang, and K. Moon, "Fast horizon detection in maritime images using region-of-interest," *Int. J. Distrib. Sensor Netw.*, vol. 14, no. 7, pp. 1–11, 2018.
- [32] D. Yongshou, L. Bowen, L. Ligang, J. Jiucui, S. Weifeng, and S. Feng, "Sea-sky-line detection based on local Otsu segmentation and Hough transform," *Opto-Electron. Eng.*, vol. 45, no. 7, 2018, Art. no. 180039.
- [33] Y. Sun and L. Fu, "Coarse-fine-stitched: A robust maritime horizon line detection method for unmanned surface vehicle applications," *Sensors*, vol. 18, no. 9, p. 2825, Aug. 2018.
- [34] T. Praczyk, "A quick algorithm for horizon line detection in marine images," *J. Mar. Sci. Technol.*, vol. 23, no. 1, pp. 164–177, Mar. 2018.
- [35] A. Carrio, H. Bavle, and P. Campoy, "Attitude estimation using horizon detection in thermal images," *Int. J. Micro Air Vehicles*, vol. 10, no. 4, pp. 352–361, Dec. 2018.
- [36] G. Ligorio and A. M. Sabatini, "Extended Kalman filter-based methods for pose estimation using visual, inertial and magnetic sensors: Comparative analysis and performance evaluation," *Sensors*, vol. 13, no. 2, pp. 1919–1941, 2013.
- [37] H.-J. Chien, C.-C. Chuang, C.-Y. Chen, and R. Klette, "When to use what feature? SIFT, SURF, ORB, or A-KAZE features for monocular visual odometry," in *Proc. Int. Conf. Image Vis. Comput. New Zealand (IVCNZ)*, Nov. 2016, pp. 1–6.
- [38] E. Hong and J. Lim, "Visual-inertial odometry with robust initialization and online scale estimation," *Sensors*, vol. 18, no. 12, p. 4287, Dec. 2018.
- [39] J. Zhang, V. Ila, and L. Kneip, "Robust visual odometry in underwater environment," in *Proc. MTS/IEEE Kobe Techno-Oceans (OTO)*, May 2018, pp. 1–9.
- [40] M. Ferrera, J. Moras, P. Trouvé-Peloux, and V. Creuze, "Real-time monocular visual odometry for turbid and dynamic underwater environments," *Sensors*, vol. 19, no. 3, p. 687, Feb. 2019.
- [41] S. Rahman, A. Q. Li, and I. Rekleitis, "SVIn2: An underwater SLAM system using sonar, visual, inertial, and depth sensor," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst. (IROS)*, Nov. 2019, pp. 1861–1868.
- [42] M. Quan, S. Piao, M. Tan, and S.-S. Huang, "Accurate monocular visual-inertial SLAM using a map-assisted EKF approach," *IEEE Access*, vol. 7, pp. 34289–34300, 2019.
- [43] C. Campos, R. Elvira, J. J. G. Rodríguez, J. M. M. Montiel, and J. D. Tardos, "ORB-SLAM3: An accurate open-source library for visual, visual-inertial, and multimap SLAM," *IEEE Trans. Robot.*, vol. 37, no. 6, pp. 1874–1890, Dec. 2021.
- [44] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected convolutional networks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jul. 2017, pp. 4700–4708.
- [45] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask R-CNN," in *Proc. IEEE Int. Conf. Comput. Vis.*, Oct. 2017, pp. 2961–2969.
- [46] P. Russo, T. Tommasi, and B. Caputo, "Towards multi-source adaptive semantic segmentation," in *Proc. Int. Conf. Image Anal. Process.* Cham, Switzerland: Springer, 2019, pp. 292–301.
- [47] P. Russo, F. Di Ciaccio, and S. Troisi, "DANAE++: A smart approach for denoising underwater attitude estimation," *Sensors*, vol. 21, no. 4, p. 1526, Feb. 2021.
- [48] X. Wang, K. Yu, S. Wu, J. Gu, Y. Liu, C. Dong, Y. Qiao, and C. C. Loy, "ESRGAN: Enhanced super-resolution generative adversarial networks," in *Proc. Eur. Conf. Comput. Vis. (ECCV) Workshops*, 2018, pp. 1–16.
- [49] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proc. IEEE Int. Conf. Comput. Vis.*, vol. 2, Sep. 1999, pp. 1150–1157.
- [50] H. Bay, T. Tuytelaars, and L. Van Gool, "Surf: Speeded up robust features," in *Proc. Eur. Conf. Comput. Vis.* Berlin, Germany: Springer, 2006, pp. 404–417.
- [51] J. R. Rambach, A. Tewari, A. Pagani, and D. Stricker, "Learning to fuse: A deep learning approach to visual-inertial camera pose estimation," in *Proc. IEEE Int. Symp. Mixed Augmented Reality (ISMAR)*, Sep. 2016, pp. 71–76.

- [52] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The KITTI dataset," *Int. J. Robot. Res.*, vol. 32, no. 11, pp. 1231–1237, 2013.
- [53] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, and R. Siegwart, "The EuRoC micro aerial vehicle datasets," *Int. J. Robot. Res.*, vol. 35, no. 10, pp. 1157–1163, Sep. 2016. [Online]. Available: <http://ijr.sagepub.com/content/early/2016/01/21/0278364915620033.abstract>
- [54] C. Li, S. Wang, Y. Zhuang, and F. Yan, "Deep sensor fusion between 2D laser scanner and IMU for mobile robot localization," *IEEE Sensors J.*, vol. 21, no. 6, pp. 8501–8509, Mar. 2021.
- [55] R. Li, S. Wang, Z. Long, and D. Gu, "UnDeepVO: Monocular visual odometry through unsupervised deep learning," in *Proc. IEEE Int. Conf. Robot. Automat. (ICRA)*, May 2018, pp. 7286–7291.
- [56] Y. Almalioglu, M. R. U. Saputra, P. P. B. D. Gusmao, A. Markham, and N. Trigoni, "GANVO: Unsupervised deep monocular visual odometry and depth estimation with generative adversarial networks," in *Proc. Int. Conf. Robot. Automat. (ICRA)*, May 2019, pp. 5474–5480.
- [57] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 27, 2014, pp. 1–9.
- [58] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele, "The cityscapes dataset for semantic urban scene understanding," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 3213–3223.
- [59] T. Feng and D. Gu, "SGANVO: Unsupervised deep visual odometry and depth estimation with stacked generative adversarial networks," *IEEE Robot. Autom. Lett.*, vol. 4, no. 4, pp. 4431–4437, Oct. 2019.
- [60] C. Zhao, Q. Sun, C. Zhang, Y. Tang, and F. Qian, "Monocular depth estimation based on deep learning: An overview," *Sci. China Technol. Sci.*, vol. 63, pp. 1–16, Jun. 2020.
- [61] P. Bernal-Polo and H. M. Barberá, "Orientation estimation by means of extended Kalman filter, quaternions, and charts," *J. Phys. Agents*, vol. 8, no. 1, pp. 1–14, 2017.
- [62] Android. *Sensorsmanager.java*. Accessed: Nov. 6, 2021. [Online]. Available: https://developer.android.com/guide/topics/sensors/sensors_overview
- [63] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2014, *arXiv:1409.1556*.
- [64] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2016, pp. 770–778.
- [65] Y. He, C. Zhu, J. Wang, M. Savvides, and X. Zhang, "Bounding box regression with uncertainty for accurate object detection," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2019, pp. 2888–2897.
- [66] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, Jun. 2015, pp. 3431–3440.
- [67] A. Veit, M. J. Wilber, and S. Belongie, "Residual networks behave like ensembles of relatively shallow networks," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 29, May 2016, pp. 550–558.
- [68] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A large-scale hierarchical image database," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, Jun. 2009, pp. 248–255.
- [69] N. Ruiz, E. Chong, and J. M. Rehg, "Fine-grained head pose estimation without keypoints," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, Jun. 2018, pp. 2074–2083.
- [70] Bosh. *BMI260: IMU Combining Accelerometer and Gyroscope*. Accessed: Oct. 1, 2021. [Online]. Available: <https://www.bosch-sensortec.com/products/motion-sensors/imus/bmi260/>
- [71] MEMSIC. *Monolithic, High Performance, Low Cost 3-Axis Magnetic Sensor*. Accessed: Oct. 1, 2021. [Online]. Available: <http://www.memsic.com/uploadfiles/2020/08/20200827165137254.pdf>
- [72] OnePlus. *OnePlus Nord—Specs*. Accessed: Oct. 1, 2021. [Online]. Available: <https://www.oneplus.com/U.K./nord-specs>
- [73] Sony. *Sony Releases Stacked CMOS Image Sensor for Smartphones*. Accessed: Oct. 1, 2021. [Online]. Available: <https://www.sony.com/en/SonyInfo/News/Press/201807/18-060E/>
- [74] AnywhereSoftware. *Simple, Powerful and Modern Development Tools*. Accessed: Oct. 4, 2021. [Online]. Available: <https://www.b4x.com/>
- [75] M. Podpora, G. P. Korbass, and A. Kawala-Janik, "YUV vs RGB-choosing a color space for human-machine interaction," in *Proc. FedCSIS*, 2014, pp. 29–34.
- [76] OpenSourceProject. *Sensors Overview*. Accessed: Nov. 6, 2021. [Online]. Available: https://github.com/aosp-mirror/platform_frameworks_base/blob/master/core/java/android/hardware/SensorManager.java
- [77] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, *arXiv:1412.6980*.
- [78] N. Otsu, "A threshold selection method from gray-level histograms," *IEEE Trans. Syst., Man, Cybern.*, vol. SMC-9, no. 1, pp. 62–66, Jan. 1979.
- [79] S. M. Ettinger, M. C. Nechyba, P. G. Ifju, and M. Waszak, "Vision-guided flight stability and control for micro air vehicles," *Adv. Robot.*, vol. 17, no. 7, pp. 617–640, 2003.
- [80] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "MobileNets: Efficient convolutional neural networks for mobile vision applications," 2017, *arXiv:1704.04861*.
- [81] S. Mittal, "A survey on optimized implementation of deep learning models on the NVIDIA Jetson platform," *J. Syst. Archit.*, vol. 97, pp. 428–442, Jan. 2019.
- [82] HUAWEI Device Co. (2021). *HUAWEI P9*. [Online]. Available: <https://consumer.huawei.com/uk/support/phones/p9/>



FABIANA DI CIACCIO received the B.S. degree in nautical and aeronautical sciences and the master's degree (*cum laude*) in science and technology of the navigation from the Parthenope University of Naples, Italy, in 2015 and 2018, respectively, where she is currently pursuing the Ph.D. degree with the International Program "Environment, Resources and Sustainable Development."

She is currently working to develop new solutions for attitude estimation based on the use of deep learning methods. Her research interests include navigation and positioning, attitude estimation, deep learning and computer vision, 3D modeling, and geomatics sciences.



PAOLO RUSSO received the B.S. degree in telecommunication engineering from the Università degli studi di Cassino, Italy, in 2008, and the M.S. degree in artificial intelligence and robotics and the Ph.D. degree in computer science from the University of Rome La Sapienza, Italy, in 2016 and 2020, respectively

From 2018 to 2019, he has been a Researcher with the Italian Institute of Technology (IIT), Tourin, Italy. He is currently an Assistant Researcher with the Alcor Laboratory, DIAG Department, University of Rome Sapienza. His main research interests include deep learning, computer vision, generative adversarial networks, and reinforcement learning.



SALVATORE TROISI received the degree (Hons.) in nautical sciences from the Faculty of Nautical Sciences, Naval University in Naples, with an experimental thesis in nautical astronomy on the use of optical amplification of light.

Since 1987, he has been a Researcher with the Group 135 (first discipline complements topography), Faculty of Nautical Sciences, Naval University of Naples. From November 1998 to September 2007, he served as an Associate Professor of SSD ICAR 06 with the Faculty of Nautical Sciences, Naval University of Naples. He has been a Full Professor SSD ICAR06 with the Faculty of Science and Technology, Parthenope University of Naples, since 2007. His research interests include deformations control networks, geoid by astrogeodetic methods and GPS, topographic methods in environmental emergencies, GPS survey for deformations, design and simulation of geodetic networks by GPS methodology, design of satellite constellations, laser scanning, filtering of laser scanning data, close-range photogrammetry for reverse engineering, and 3D building modeling by aerial laser scanning data.