

Toward a Generic Diver-Following Algorithm: Balancing Robustness and Efficiency in Deep Visual Detection

Md. Jahidul Islam^{1b}, Michael Fulton, and Junaed Sattar^{1b}

Abstract—This letter explores the design and development of a class of robust diver detection algorithms for autonomous diver-following applications. By considering the operational challenges for underwater visual tracking in diverse real-world settings, we formulate a set of desired features of a generic diver-following algorithm. We attempt to accommodate these features and maximize general tracking performance by exploiting the state-of-the-art deep object detection models. We fine tune the building blocks of these models with a goal of balancing the tradeoff between robustness and efficiency in an on-board setting under real-time constraints. Subsequently, we design an architecturally simple convolutional neural network based diver detection model that is much faster than the state-of-the-art deep models yet provides comparable detection performances. In addition, we validate the performance and effectiveness of the proposed model through a number of diver-following experiments in closed-water and open-water environments.

Index Terms—Human detection and tracking, field robots, marine robotics.

I. INTRODUCTION

UNDERWATER applications of autonomous underwater robots range from inspection and surveillance to data collection and mapping tasks. Such missions often require a team of divers and robots to collaborate for successful completion. Without sacrificing the generality of such applications, we can consider a single-robot setting where a human diver leads the task and interacts with the robot which follows the diver at certain stages of the mission. Such situations arise in numerous important applications such as submarine pipeline and shipwreck inspection, marine life and seabed monitoring, and many other underwater exploration activities [1]. Although following the diver is not the primary objective in these applications, it significantly simplifies the operational loop and reduces the associated overhead by eliminating the necessity of tele-operation.

Robust underwater visual perception is generally challenging due to marine artifacts such as poor visibility, variations in

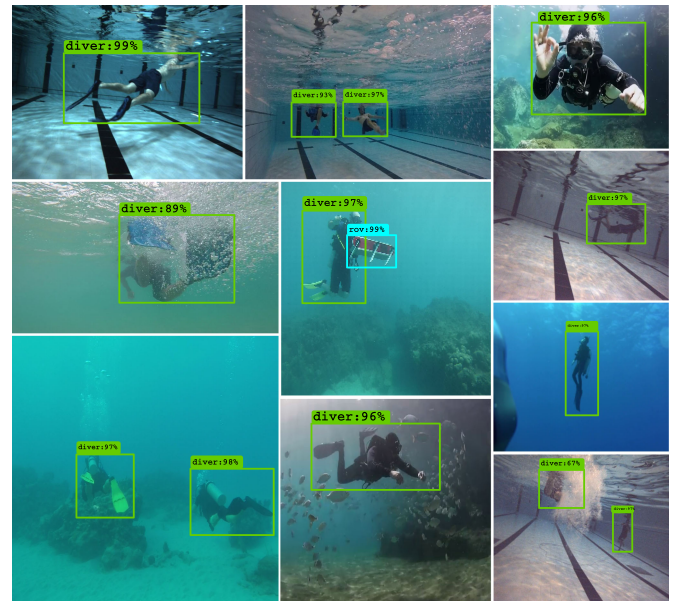


Fig. 1. Snapshots of a set of diverse first-person views of the robot from different diver-following scenarios. Notice the variation in appearances of the divers and possible noise or disturbances in the scene over different scenarios. The rectangles and text overlaid on the figures are the outputs generated by our model at test time.

illumination, suspended particles, etc. Additionally, color distortion and scarcity of salient visual features make it harder to robustly detect and accurately follow a diver in arbitrary directions. Moreover, divers' appearances to the robot vary greatly based on their swimming styles, choices of wearables, and relative orientations with respect to the robot. These problems are exacerbated underwater since both the robot and diver are suspended in a six-degrees-of-freedom (6DOF) environment. Consequently, classical model-based detection algorithms fail to achieve good generalization performance [2]. On the other hand, model-free algorithms incur significant target drift [3] under such noisy conditions.

In this letter, we address the inherent difficulties of underwater visual detection by designing a class of diver detection algorithms that are: a) invariant to color (of divers' body/wearables [4]), b) invariant to divers' relative motion and orientation, c) robust to noise and image distortions [5], and d) reasonably efficient for real-time deployment. We exploit the current state-of-the-art object detectors to accommodate these features and

Manuscript received September 11, 2018; accepted November 5, 2018. Date of publication November 22, 2018; date of current version November 30, 2018. This letter was recommended for publication by Associate Editor E. Johnson and Editor J. Roberts upon evaluation of the reviewers' comments. (Corresponding author: Md. Jahidul Islam.)

The authors are with the Interactive Robotics and Vision Laboratory, Department of Computer Science and Engineering, University of Minnesota, Twin Cities, Minneapolis MN 55455 USA (e-mail: islam034@umn.edu; fulto081@umn.edu; junaed@umn.edu).

Digital Object Identifier 10.1109/LRA.2018.2882856

maximize the generalization performance for diver detection using RGB images as input. Specifically, we use the following four models: Faster R-CNN [6] with Inception V2 [7] as a feature extractor, Single Shot MultiBox Detector (SSD) [8] with MobileNet V2 [9], [10] as a feature extractor, You Only Look Once (YOLO) V2 [11], and Tiny YOLO [12]. These are the fastest (in terms of processing time of a single frame) among the family of current state-of-the-art models [13] for general object detection. We train these models using a rigorously prepared dataset containing sufficient training instances to capture the variabilities of underwater visual sensing.

Subsequently, we design an architecturally simple (i.e., sparse) CNN-based model that is computationally much faster than the state-of-the-art diver detection models. The faster running time ensures real-time tracking performance with limited on-board computational resources. We also demonstrate its effectiveness in terms of detection performances compared to the state-of-the-art models through extensive quantitative experiments. We then validate these results with a series of field experiments. Based on our design, implementation, and experimental findings, we make the following contributions in this letter:

- We attempt to overcome the limitations of existing model-based diver-following algorithms by leveraging state-of-the-art deep object detection models. These models are trained on comprehensive datasets to deal with the challenges involved in underwater visual perception.¹
- In addition, we design a CNN-based diver detection model to balance the trade-offs between robustness and efficiency. The proposed model provides considerably faster running time, in addition to achieving detection performances comparable to the state-of-the-art models.
- Finally, we validate the effectiveness of the proposed diver detection models through extensive experimental evaluations. A number of diver-following experiments are performed both in open-water and closed-water (i.e., oceans and pools, respectively) environments in order to demonstrate their real-time tracking performances.

Furthermore, we demonstrate that the proposed models can be extended for a wide range of other applications such as human-robot communication [14], robot convoying [3], cooperative localization [15], [16], etc. The state-of-the-art detection performance, fast running time, and architectural portability are the key features of these models, which make them suitable for underwater human-robot collaborative applications.

II. RELATED WORK

A categorization of the visual perception techniques that are commonly used for autonomous diver-following is illustrated in Fig. 2. Based on algorithmic usage of the input features, the perception techniques can be classified as feature-based tracking, feature-based learning, or feature/representation learning algorithms. On the other hand, they can be categorized as model-based or model-free techniques based on whether or not any

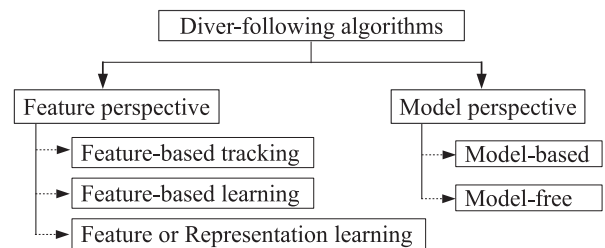


Fig. 2. An algorithmic categorization of the visual perception techniques used for diver-following [17].

prior knowledge about the appearance or motion of the diver is used for tracking.

A. Model Perspective

In model-free algorithms, no prior information about the target (e.g., diver's motion model, color of wearables, etc.) is used for tracking. These algorithms are initialized arbitrarily and then iteratively learn to track the target in a semi-supervised fashion [18]. TLD ("tracking-learning-detection") trackers [19] and optical flow-based trackers [20] are the most commonly used model-free algorithms for general object tracking. The TLD trackers train a detector using positive and negative feedback that are obtained from image-based features. In contrast, the optical flow-based methods estimate the motion of each pixel by solving the Horn and Schunck formulation [21]. Although model-free techniques work reasonably well in practice for general object tracking, they often suffer from tracking drift caused by the accumulation of detection errors over time.

On the other hand, model-based algorithms use prior knowledge about the divers' motion and appearances in order to formulate a model in the input feature-space. Iterative search methods are then applied to find the target model in the feature-space [22]. Machine learning techniques are also widely used to learn the diver-specific features [17], [23] and predict the target location in the feature-space. Performance of the model-free algorithms depend on comprehensiveness of the model descriptors and the underlying input feature-space. Hence, they require careful design and thorough training processes to ensure good tracking performance.

B. Feature Perspective

Simple feature-based trackers [24], [25] are often practical choices for autonomous diver-following due to their operational simplicity and computational efficiency. For instance, color-based trackers perform binary image thresholding based on the color of a diver's flippers or suit. The thresholded binary image is then refined to track the centroid of the target (diver) using algorithms such as mean-shift, particle filters, etc. Optical flow-based methods can also be utilized to track divers' motion in the spatio-temporal volume [17], [21].

Since color distortions and low visibility issues are common in underwater settings, frequency-domain signatures of divers' swimming patterns are often used for reliable detection. Specifically, intensity variations in the spatio-temporal volume caused

¹The dataset and trained models are available for academic research purposes.

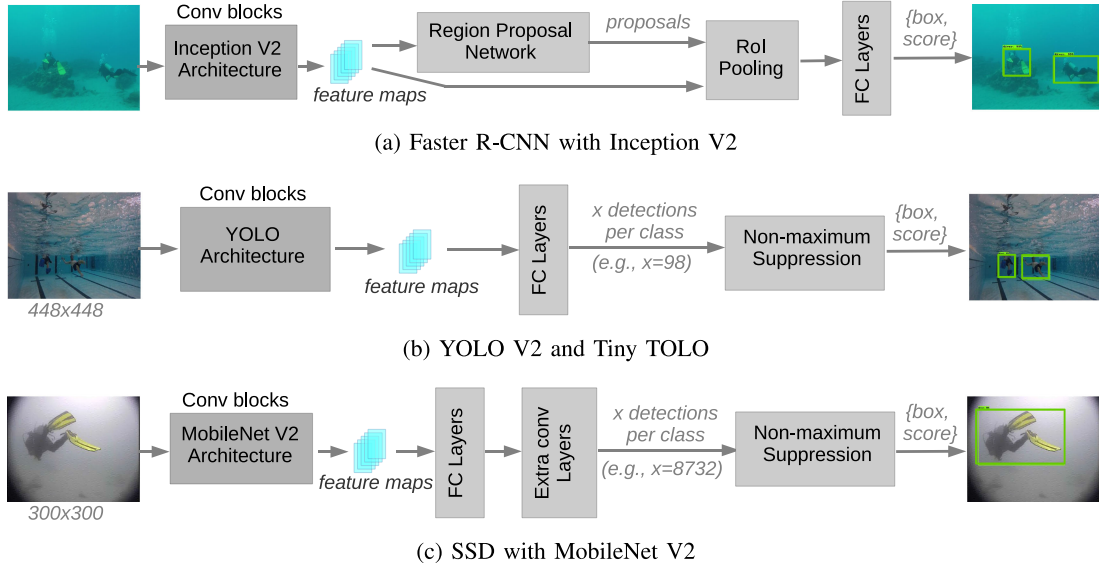


Fig. 3. Schematic diagrams of the deep visual models used for diver detection.

by a diver's swimming gait generate identifiable high-energy responses in the 1–2 Hz frequency range, which can be used for diver detection [26]. Moreover, the frequency-domain signatures can be combined with the spatial-domain features for robust diver tracking. For instance, in [22], a Hidden Markov Model (HMM) is used to track divers' potential swimming trajectories in the spatio-temporal domain, and then frequency-domain features are utilized to detect the diver along those trajectories.

Another class of approaches use machine learning techniques to approximate the underlying function that relates the input feature-space to the target model of the diver. For instance, Support Vector Machines (SVMs) are trained using Histogram of Oriented Gradients (HOG) features [27] for robust person detection in general. Ensemble methods such as Adaptive Boosting (AdaBoost) [23] are also widely used as they are computationally inexpensive yet highly accurate in practice. AdaBoost learns a strong tracker from a large number of simple feature-based diver trackers. Several other machine learning techniques have been investigated for diver tracking and underwater object tracking in general [17]. One major challenge involved in using these models is to design a set of robust features that are invariant to noise, lighting condition, and other variabilities such as divers' swimming motion and wearables.

Convolutional Neural Network (CNN)-based deep models improve generalization performance by learning a feature representation from the image-space. The extracted features are used as inputs to the detector (i.e., fully-connected layers); this end-to-end training process significantly improves the detection performance compared to using hand-crafted features. Once trained with sufficient data, these models are quite robust to occlusion, noise, and color distortions [3]. Despite the robust performance, the applicability of these models to real-time applications is often limited due to their slow running time on embedded devices. In this letter, we investigate the performances and feasibilities of the state-of-the-art deep object detectors for diver-following

applications. We also design a CNN-based model that achieves robust detection performance in addition to ensuring that the real-time operating constraints are met.

III. NETWORK ARCHITECTURE AND DESIGN

A. State-of-the-Art Object Detectors

We use a Faster R-CNN model, two YOLO models, and an SSD model for diver detection. These are end-to-end trainable models and provide state-of-the-art performances on standard object detection datasets; we refer to [12], [13] for detailed comparisons of their detection performances and running times. As outlined in Fig. 3, we now briefly discuss their methodologies and the related design choices in terms of major computational components.

1) *Faster R-CNN With Inception V2*: Faster R-CNN [6] is an improvement of R-CNN [28] that introduces a Region Proposal Network (RPN) to make the whole object detection network end-to-end trainable. The RPN uses the last convolutional feature-maps to produce region proposals which are then fed to the fully connected layers for the final detection. The original implementation of Faster R-CNN uses VGG-16 [29] model for feature extraction. However, we use Inception V2 [7] model for feature extraction instead, as it is known to provide better object detection performances on standard datasets [13].

2) *YOLO V2 and Tiny YOLO*: YOLO models [11], [30] formulate object detection as a regression problem in order to avoid using computationally expensive RPNs. They divide the image-space into rectangular grids and predict a fixed number of bounding boxes, their corresponding confidence scores, and class probabilities. Although there are restrictions on the maximum number of object categories, they perform faster than the standard RPN-based object detectors. Tiny YOLO [12] is a scaled down version of the original model having sparser layers that runs much faster compared to the original model; however, it sacrifices detection accuracy in the process.

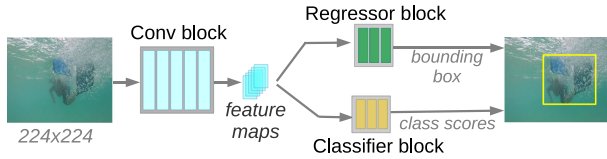


Fig. 4. A schematic diagram of the proposed CNN-based model for detecting a single diver in the image-space.

TABLE I

PARAMETERS AND DIMENSIONS OF THE CNN MODEL OUTLINED IN FIG. 4. (CONVOLUTIONAL BLOCK: CONV1-CONV5, CLASSIFIER BLOCK: FC1-FC3, REGRESSION BLOCK: RC1-RC3; N: THE NUMBER OF OBJECT CATEGORIES; *AN ADDITIONAL POOLING LAYER WAS USED BEFORE PASSING THE CONV5 FEATURES-MAPS TO FC1)

Layer	Input feature-map	Kernel size	Strides	Output feature-map
conv1	224x224x3	11x11x3x64	[1,4,4,1]	56x56x64
pool1	56x56x64	1x3x3x1	[1,2,2,1]	27x27x64
conv2	27x27x64	5x5x64x192	[1,1,1,1]	27x27x192
pool2	27x27x192	1x3x3x1	[1,2,2,1]	13x13x192
conv3	13x13x192	3x3x192x192	[1,1,1,1]	13x13x192
conv4	13x13x192	3x3x192x192	[1,1,1,1]	13x13x192
conv5	13x13x192	3x3x192x128	[1,1,1,1]	13x13x128
fc1	4608x1*	—	—	1024x1
fc2	1024x1	—	—	128x1
fc3	128x1	—	—	n
rc1	21632x1	—	—	4096x1
rc2	4096x1	—	—	192x1
rc3	192x1	—	—	4n

3) *SSD With MobileNet V2*: SSD (Single-Shot Detector) [8] also performs object localization and classification in a single pass of the network using the regression trick as in the YOLO [30] model. The architectural difference of SSD with YOLO is that it introduces additional convolutional layers to the end of a base network, which results in improved performances. In our implementation, we use MobileNet V2 [9] as the base network to ensure faster running time.

B. Proposed CNN-Based Model

Fig. 4 shows a schematic diagram of the proposed CNN-based diver detection model. It consists of three major parts: a convolutional block, a regressor block, and a classifier block. The convolutional block consists of five layers, whereas the classifier and regressor block each consist of three fully connected layers. Detailed network parameters and dimensions are specified in Table I.

1) *Design Intuition*: The state-of-the-art deep visual models are designed for general applications and are trained on standard datasets having a large number of object categories. However, for most underwater human-robot collaborative applications including diver-following, only a few object categories (e.g., diver, robot, coral reefs, etc.) are relevant. We try to take advantage of this by designing an architecturally simpler model that ensures much faster running time in an embedded platform in addition to providing robust detection performance. The underlying design intuitions can be summarized as follows:

- The proposed model demonstrated in Fig. 4 is particularly designed for detecting a single diver. Five convolutional

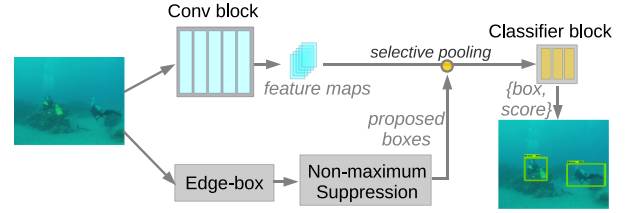


Fig. 5. Allowing detections of multiple divers in the proposed model using a region selector named Edge-box [31].

layers are used to extract the spatial features in the RGB image-space by learning a set of convolutional kernels.

- The extracted features are then fed to the classifier and regressor block for detecting a diver and localizing the corresponding bounding box, respectively. Both the classifier and regressor block consist of three fully connected layers.
- Therefore, the task of the regressor block is to locate a potential diver in the image-space, whereas the classifier block provides the confidence scores associated with that detection.

The proposed model has a sparse convolutional block and uses a three layer regressor block instead of using an RPN. As demonstrated in Table I, it has significantly fewer network parameters compared to the state-of-the-art object detection models.

2) *Allowing Multiple Detections*: Although following a single diver is the most common diver-following scenario, detecting multiple divers and other objects is necessary for many human-robot collaborative applications. As shown in Fig. 5, we add multi-object detection capabilities in our proposed model by replacing the regressor with a region selector. We use the state-of-the-art class-agnostic region selector named Edge-box [31]. Edge-box utilizes the image-level statistics like edges and contours in order to measure *objectness scores* in various prospective regions in the image-space.

We use the same convolutional block to extract feature maps. The bounding boxes generated by Edge-box are filtered based on their objectness scores and then non-maxima suppression techniques are applied to get the dominant regions of interest in the image-space. The corresponding feature maps are then fed to the classifier block to predict the object categories. Although we need additional computation for Edge-box, it runs independently and in parallel with the convolutional block; the overall pipeline is still faster than if we were to use an RPN-based object detector model.

IV. EXPERIMENTS

We now discuss the implementation details of the proposed networks and present the experimental results.

A. Dataset Preparation

We performed numerous diver-following experiments in pools and oceans in order to prepare training datasets for the deep models. In addition, we collected data from underwater

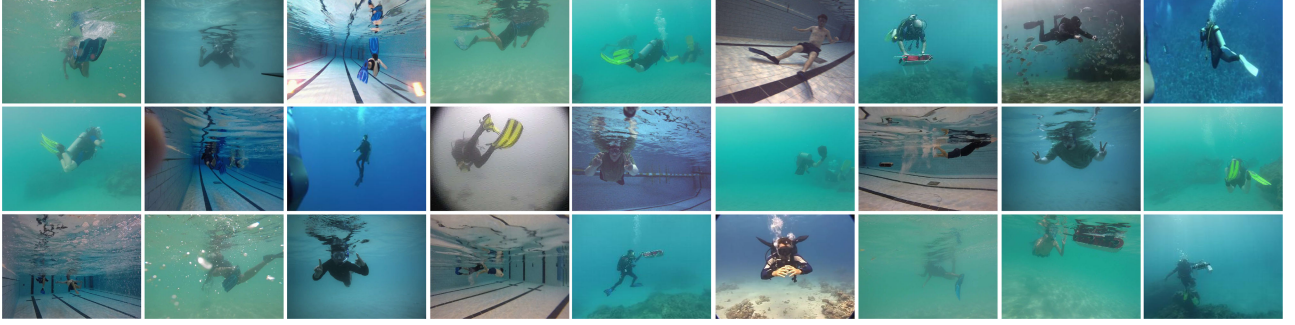


Fig. 6. A few samples from the training dataset are shown. The annotated training images have class labels (e.g., diver, robot) and corresponding bounding boxes. A total of 30 K of these annotated images are used for supervised training.

field trials that are performed by different research groups over the years in pools, lakes, and oceans. This variety of experimental setups is crucial to ensure comprehensiveness of the datasets so that the supervised models can learn the inherent diversity of various application scenarios. We made sure that the datasets contain training examples to capture the following variabilities:

- *Natural variabilities*: changes in visibilities for different sources of water, lighting conditions at varied depths, chromatic distortions, etc.
- *Artificial variabilities*: data collected using different robots and cameras.
- *Human variabilities*: different persons and appearances, choice and variations of wearables such as suits, flippers, goggles, etc.

We extracted the robot's camera-feed during these experiments and prepared image-based datasets for supervised training. The images are annotated using the 'label-image' software (github.com/tzutalin/labelImg) by a number of human participants (acknowledged later in the letter) over the period of six months. Few sample images from the dataset are shown in Fig. 6; it contains a total of 30 K images, which are annotated to have class-labels and bounding boxes.

B. Supervised Training Processes

We train all the supervised deep models on a Linux machine with four GPU cards (NVIDIA GTX 1080). TensorFlow [32] and Darknet [12] libraries are used for implementation. Once the training is done, the trained inference model (and parameters) is saved and transferred to the robot CPU for validation and real-time experiments.

For the state-of-the-art models (Fig. 3), we utilized the pre-trained models for Faster R-CNN, YOLO, and SSD. These models are trained with the recommended configurations provided with their APIs; we refer to [12], [13] for the detailed processes. On the other hand, our proposed CNN-based models are trained from scratch. Non-supervised pre-training and drop-outs are not used while training. RMSProp [33] is used as the optimization function with an initial learning rate of 0.001. In addition, standard cross-entropy and L_2 loss functions are used by the classifier and regressor, respectively. Visualization for the overall convergence behavior of the model is provided in Fig. 7.

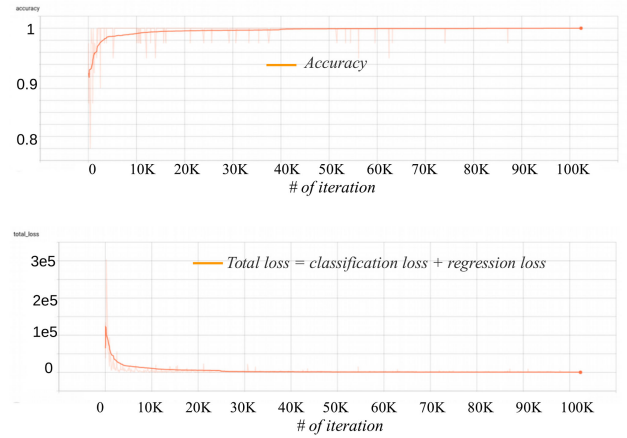


Fig. 7. Convergence behavior of the proposed CNN-based model in terms of training accuracy (top) and training loss (bottom).

C. Performance Evaluation

We evaluate and compare detection performances of all the models based on standard performance metrics. The test dataset contain 2.2 K images that are chosen from separate field experiments (i.e., they are excluded from the training dataset).

1) *Metrics*: We use the following two standard performance metrics:

- *mAP (mean Average Precision)*: it is the average of the maximum precisions at different recall values. The precision and recall are defined as $precision = \frac{TP}{TP+FP}$ and $recall = \frac{TP}{TP+FN}$; here, the terms TP, FP, and FN are short forms of True Positive, False Positive, and False Negative, respectively.
- *IoU (Intersection over Union)*: it is a measure of how well a model predicts the locations of the objects. It is calculated using the area of overlapping regions of the predicted and ground truth bounding boxes, defined as $IoU = \frac{Area\ of\ overlap}{Area\ of\ union}$.

As their definitions suggest, mAP measures the detection accuracy, and IoU measures the object localization performance. We also evaluate and compare the running times of the models based on FPS (Frames Per Second), the (average) number of image-frames that a model can process per second. We measure the running times on three different devices:

- NVIDIA GTX 1080 GPU

TABLE II
PERFORMANCE COMPARISON FOR THE DIVER DETECTION MODELS BASED ON
STANDARD METRICS

Models	mAP (%)	IoU (%)	FPS		
			GTX 1080	Jetson TX2	Robot CPU
Faster R-CNN (Inception V2)	71.1	78.3	17.3	2.1	0.52
YOLO V2	57.84	62.42	73.3	6.2	0.11
Tiny YOLO	52.33	59.94	220	20	5.5
SSD (MobileNet V2)	61.25	69.8	92	9.85	3.8
Proposed CNN-based Model	53.75	67.4	263.5	17.35	6.85

- Embedded GPU (NVIDIA Jetson TX2)
- Robot CPU (Intel i3-6100U)

2) *Results:* The performances of the diver detection models based on mAP, IoU, and FPS are illustrated in Table II. The Faster R-CNN (Inception V2) model achieves much better detection performances compared to the other models although it is the slowest in terms of running time. On the other hand, YOLO V2, SSD (MobileNet V2), and the proposed CNN-based model provide comparable detection performances. Although Tiny YOLO provides fast running time, its detection performance is not as good as the other models.

As the results demonstrate, the proposed CNN-based model runs at a rate of 6.85 FPS on the robot CPU and 17.35 FPS on the embedded GPU, which validate its applicability for real-time diver-following applications. This fast running time comes at a cost of losing approximately 18% mAP and 11% IoU compared to the Faster R-CNN (Inception V2) model. Nevertheless, in our real-world experiments, we have found these detection performances to be sufficient for achieving reasonable tracking performances. In the following sections, we provide details of these field experiments and discuss general applicability of the proposed model from a practical standpoint.

D. Field Experiments

1) *Setup:* We have performed several real-world experiments both in closed-water and in open-water conditions (i.e., in pools and in oceans). An autonomous underwater robot of the Aqua [34] family is used for testing the diver-following modules. During the experiments, a diver swims in front of the robot in arbitrary directions. The task of the robot is to visually detect the diver using its camera feed and follow behind him/her with a smooth motion.

2) *Visual Servoing Controller:* The Aqua robots have five degrees-of-freedom of control, i.e., three angular (yaw, pitch, and roll) and two linear (forward and vertical speed) controls. In our experiments for autonomous diver-following, we adopt a tracking-by-detection method where the visual servoing [35] controller uses the uncalibrated camera feeds for navigation. The controller regulates the motion of the robot in order to bring the observed bounding box of the target diver to the center of the camera image. The distance of the diver is approximated by the size of the bounding box and forward velocity

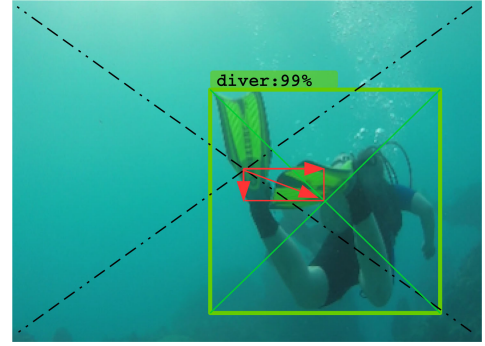


Fig. 8. Illustration of how the yaw and pitch commands are generated based on the horizontal and vertical displacements of the center of the detected bounding box.

rates are generated accordingly. Additionally, the yaw and pitch commands are normalized based on the horizontal and vertical displacements of the observed bounding box-center from the image-center (see Fig. 8); these navigation commands are then regulated by separate PID controllers. On the other hand, the roll stabilization and hovering are handled by the robot's autopilot module [36].

3) *Feasibility and General Applicability:* As mentioned, the diver-following module uses a monocular camera feed of the robot in order to detect a diver in the image-space and generate a bounding box. The visual servoing controller uses this bounding box and regulates robot motion commands in order to follow the diver. Therefore, correct detection of the diver is essential for overall success of the operation. We provided the detection performances of our proposed model over a variety of test scenarios in Table II (few snapshots are illustrated in Fig. 1). During the field experiments, we have found 6–7 positive detections per second on an average, which is sufficient for successfully following a diver in real-time [2]. In addition, the on-board memory overhead is low as the saved inference model is only about 60 MB in size.

In addition, the proposed model is considerably robust to occlusion and noise, in addition to being invariant to divers' appearances and wearables. Nevertheless, the detection performances might be negatively affected by unfavorable visual conditions; we demonstrate few such cases in Fig. 9. In Fig. 9(a), the diver is only partially detected with low confidence (67%). This is because the flippers' motion produces a flurry of air-bubbles (since he was swimming very close to the ocean surface), which occluded the robot's view. Suspended particles cause similar difficulties in diver-following scenarios. The visual servoing controller can recover from such inaccurate detections as long as the diver is partially visible. However, the continuous tracking might fail if the diver moves away from the robot's field of view before it can recover. In this experiment, 27 consecutive inaccurate detections (i.e., confidence score less than 50%) caused enough drift in the robot's motion for it to lose sight of the person. On the other hand, occlusion also affects the detection performances as shown in Fig. 9(b); here, the proposed model could not localize the two divers correctly due to occlusion.

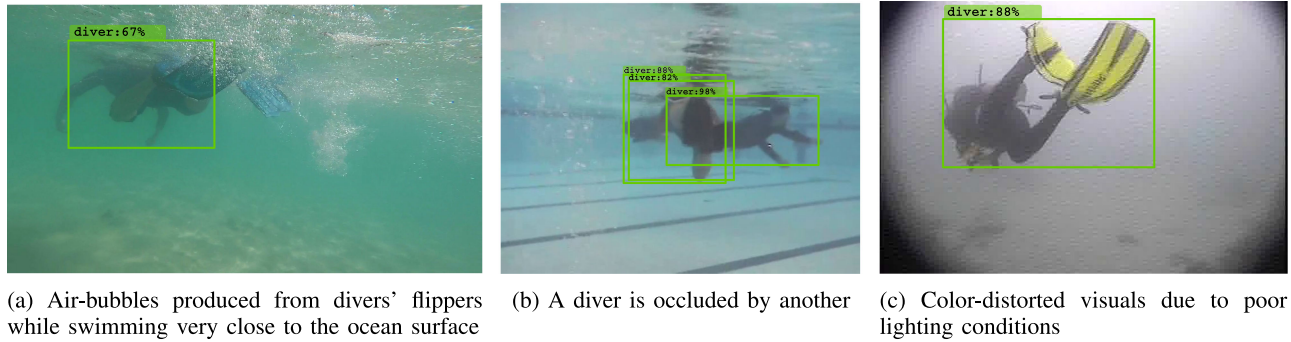


Fig. 9. A few cases where the diver-detection performance is challenged by noise and occlusion.

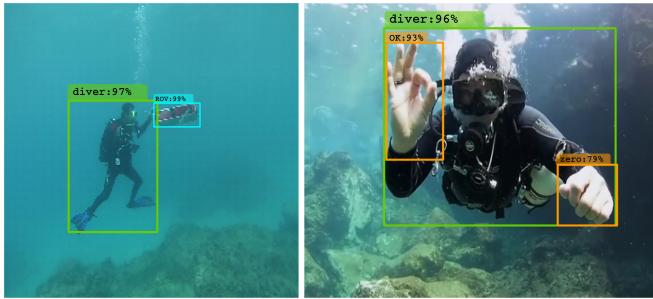


Fig. 10. Detection of ROVs and hand gestures by the same diver-detector model. In this case, the SSD (MobileNet V2) model was re-trained on additional data and object categories for ROV and hand gestures (used for human-robot communication [14]).

Lastly, since our training datasets include a large collection of gray-scale and color distorted underwater images, the proposed models are considerably robust to noise and color distortions (Fig. 9(c)). Nonetheless, state-of-the-art image enhancement techniques for underwater imagery can be utilized to alleviate severe chromatic distortions. We refer interested readers to [5], where we tried to address these issues for generic underwater applications.

We also performed experiments to explore the usabilities of the proposed diver detection models for other underwater applications. As demonstrated in Fig. 10, by simply re-training on additional data and object categories, the same models can be utilized in a wide range of underwater human-robot collaborative applications such as following a team of divers, robot conveying [3], human-robot communication [14], etc. In particular, if the application do not pose real-time constraints, we can use models such as Faster R-CNN (Inception V2) for better detection performances.

V. CONCLUSIONS

In this letter, we have tried to address the challenges involved in underwater visual perception for autonomous diver-following. At first, we investigated the performances and applicabilities of the state-of-the-art deep object detectors. We prepared and used a comprehensive dataset for training these models; then we fine-tuned each computational component in order to meet the real-time and on-board operating constraints. Subsequently, we designed a CNN-based diver detection model

that establishes a delicate balance between robust detection performance and fast running time. Finally, we validated the tracking performances and general applicabilities of the proposed models through a number of field experiments in pools and oceans.

In the future, we seek to improve the running time of the general object detection models on embedded devices. Additionally, we aim to investigate the use of human body-pose detection models to understand divers' motion, instructions, and activities.

ACKNOWLEDGMENT

The authors gratefully acknowledge the support of the Mn-Drive initiative on this research. The authors would like to thank the Bellairs Research Institute of Barbados for providing the facilities for their field experiments, the NVIDIA Corporation for its support with the donation of a Titan Xp GPU for their research, and their colleagues, namely C. Fabbri, M. Ho, E. Imhoff, Y. Xia, and J. Lagman for their assistance in collecting and annotating the training data.

REFERENCES

- [1] J. Sattar *et al.*, "Enabling autonomous capabilities in underwater robotics," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2008, pp. 3628–3634.
- [2] M. J. Islam, M. Ho, and J. Sattar, "Understanding human motion and gestures for underwater human-robot collaboration," *J. Field Robot.*, vol. 35, 2018, pp. 1–23.
- [3] F. Shkurti *et al.*, "Underwater multi-robot conveying using visual tracking by detection," in *Proc. IEEE/RSJ Intell. Conf. Intell. Robots Syst.*, 2017, pp. 4189–4196.
- [4] J. Sattar and G. Dudek, "Where is your dive buddy: Tracking humans underwater using spatio-temporal features," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2007, pp. 3654–3659.
- [5] C. Fabbri, M. J. Islam, and J. Sattar, "Enhancing underwater imagery using generative adversarial networks," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2018, pp. 7159–7165.
- [6] S. Ren, K. He, R. Girshick, and J. Sun, "Faster R-CNN: Towards real-time object detection with region proposal networks," in *Proc. Adv. Neural Inf. Process. Syst.*, 2015, pp. 91–99.
- [7] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2818–2826.
- [8] W. Liu *et al.*, "SSD: Single shot multibox detector," in *Proc. Eur. Conf. Comput. Vis.*, 2016, pp. 21–37.
- [9] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 4510–4520.
- [10] A. Howard *et al.*, "MobileNets: Efficient convolutional neural networks for mobile vision applications," *Comput. Res. Repository (CoRR)*, abs/1704.04861, pp. 1–9, 2017.

- [11] J. Redmon and A. Farhadi, "YOLO9000: Better, faster, stronger," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 6517–6525.
- [12] J. Redmon and A. Farhadi, "Tiny YOLO," 2017. [Online]. Available: <https://pjreddie.com/darknet/yolo/>. Accessed on: Oct. 10, 2018.
- [13] Google, "TensorFlow object detection zoo," 2017. [Online]. Available: https://github.com/tensorflow/models/blob/master/research/object_detection/g3doc/detection_model_zoo.md. Accessed on: Oct. 10, 2018.
- [14] M. J. Islam, M. Ho, and J. Sattar, "Dynamic reconfiguration of mission parameters in underwater human-robot collaboration," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2018, pp. 1–8.
- [15] A. Bahr, J. J. Leonard, and M. F. Fallon, "Cooperative localization for autonomous underwater vehicles," *Int. J. Robot. Res.*, vol. 28, no. 6, pp. 714–728, 2009.
- [16] I. Rekleitis, G. Dudek, and E. Milios, "Probabilistic cooperative localization and mapping in practice," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2003, vol. 2, pp. 1907–1912.
- [17] M. J. Islam, J. Hong, and J. Sattar, "Person following by autonomous robots: A categorical overview," *Comput. Res. Repository (CoRR)*, abs/1803.08202, pp. 1–25, 2018.
- [18] Q. Yu, T. B. Dinh, and G. Medioni, "Online tracking and reacquisition using co-trained generative and discriminative trackers," in *Proc. Eur. Conf. Comput. Vis.*, 2008, pp. 678–691.
- [19] Z. Kalal, K. Mikolajczyk, and J. Matas, "Tracking-learning-detection," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 34, no. 7, pp. 1409–1422, Jul. 2012.
- [20] J. Shin *et al.*, "Optical flow-based real-time object tracking using non-prior training active feature model," *Real-Time Imag.*, vol. 11, no. 3, pp. 204–218, 2005.
- [21] H. Inoue, T. Tachikawa, and M. Inaba, "Robot vision system with a correlation chip for real-time tracking, optical flow and depth map generation," in *Proc. IEEE Int. Conf. Robot. Automat.*, 1992, pp. 1621–1626.
- [22] M. J. Islam and J. Sattar, "Mixed-domain biological motion tracking for underwater human-robot interaction," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2017, pp. 4457–4464.
- [23] J. Sattar and G. Dudek, "Robust servo-control for underwater robots using banks of visual filters," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2009, pp. 3583–3588.
- [24] J. Sattar and G. Dudek, "On the performance of color tracking algorithms for underwater robots under varying lighting and visibility," in *Proc. IEEE Int. Conf. Robot. Automat.*, 2006, pp. 3550–3555.
- [25] J. Sattar, P. Giguere, G. Dudek, and C. Prahacs, "A visual servoing system for an aquatic swimming robot," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, 2005, pp. 1483–1488.
- [26] J. Sattar and G. Dudek, "Underwater human-robot interaction via biological motion identification," *Robot. Sci. Syst.*, vol. 5, pp. 185–192, 2009.
- [27] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2005, vol. 1, pp. 886–893.
- [28] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2014, pp. 580–587.
- [29] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *Comput. Res. Repository (CoRR)*, abs/1409.1556, pp. 1–14, 2014.
- [30] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 779–788.
- [31] C. L. Zitnick and P. Dollár, "Edge boxes: Locating object proposals from edges," in *Proc. Eur. Conf. Comput. Vis.*, 2014, pp. 391–405.
- [32] M. Abadi *et al.*, "Tensorflow: Large-scale machine learning on heterogeneous distributed systems," in *Proc. USENIX Conf. Operating Syst. Des. Implementation (OSDI)*, 2016, pp. 265–283.
- [33] T. Tieleman and G. Hinton, "Lecture 6.5-RMSProp: Divide the gradient by a running average of its recent magnitude," *COURSERA: Neural Netw. Mach. Learn.*, vol. 4, no. 2, pp. 26–31, 2012.
- [34] G. Dudek *et al.*, "Aqua: An amphibious autonomous robot," *Computer*, vol. 40, no. 1, pp. 46–53, Jan. 2007.
- [35] B. Espiau, F. Chaumette, and P. Rives, "A new approach to visual servoing in robotics," *IEEE Trans. Robot. Autom.*, vol. 8, no. 3, pp. 313–326, Jun. 1992.
- [36] D. Meger, F. Shkurti, D. C. Poza, P. Giguere, and G. Dudek, "3d trajectory synthesis and control for a legged swimming robot," in *Proc. IEEE/RSJ Intell. Conf. Robots Syst.*, 2014, pp. 2257–2264.