

DeepURL: Deep Pose Estimation Framework for Underwater Relative Localization

Bharat Joshi^{a*}, Md Modasshir^{a*}, Travis Manderson^b, Hunter Damron^a, Marios Xanthidis^a,
Alberto Quattrini Li^c, Ioannis Rekleitis^a, Gregory Dudek^b.

Abstract—In this paper, we propose a real-time deep learning approach for determining the 6D relative pose of Autonomous Underwater Vehicles (AUV) from a single image. A team of autonomous robots localizing themselves in a communication-constrained underwater environment is essential for many applications such as underwater exploration, mapping, multi-robot convoying, and other multi-robot tasks. Due to the profound difficulty of collecting ground truth images with accurate 6D poses underwater, this work utilizes rendered images from the Unreal Game Engine simulation for training. An image-to-image translation network is employed to bridge the gap between the rendered and the real images producing synthetic images for training. The proposed method predicts the 6D pose of an AUV from a single image as 2D image keypoints representing 8 corners of the 3D model of the AUV, and then the 6D pose in the camera coordinates is determined using RANSAC-based PnP. Experimental results in real-world underwater environments (swimming pool and ocean) with different cameras demonstrate the robustness and accuracy of the proposed technique in terms of translation error and orientation error over the state-of-the-art methods. The code is publicly available.¹

I. INTRODUCTION

The ability to localize is crucial to many robotic applications. There are several environments where keeping track of the vehicle’s position is a challenging task; particularly in GPS-denied environments with limited features. A common approach to address the localization problem is to use intra-robot measurements for improved positional accuracy – an approach termed Cooperative Localization (CL) [1]. Central to CL is the ability to estimate the relative pose between the two robots; this estimate can then be utilized to improve the absolute localization based on the global pose estimates for one of the two robots. A video overview is also available online¹.

In this paper, we propose and evaluate a deep pose estimation framework for underwater relative localization, called *DeepURL*. The primary application motivating this work is underwater exploration and mapping by a team of

^aUniversity of South Carolina, Columbia, SC, USA, 29208, {bjoshi, modasshm, mariosx, hdamron}@email.sc.edu, yiannisr@cse.sc.edu. The authors would like to acknowledge the generous support of the National Science Foundation grants (NSF 2024741, 1943205)

^bMcGill University, Montreal, Quebec, Canada, {travism, dudek}@cim.mcgill.ca

^cDartmouth College, Hanover, NH, USA, 03755 alberto.quattrini.li@dartmouth.edu

*First two authors contributed equally.

¹https://github.com/joshi-bharat/deep_underwater_localization



Fig. 1: Two Aqua2 vehicles collecting images over a reef require relative localization to efficiently cover the area.

autonomous underwater vehicles (AUVs) with a focus on shipwreck and underwater cave mapping; environments that are challenging to most existing localization methodologies (e.g., visual and visual/inertial-based systems [2], [3]). Other applications include convoying [4], environmental assessments [5], informative navigation [6], and inspections.

The proposed methodology draws from the rich object detection research and is adapted to the unique conditions of the underwater domain. Traditionally, 6D pose estimation (3D position and 3D orientation) is performed by matching feature points between 3D models and images [7]–[10]. While these methods are robust when objects are well textured, they perform poorly when objects are featureless or textureless. In the underwater domain, particles in the water generate undesired texture smoothing. Recent approaches [11]–[16] to estimate 6D poses using deep neural network perform well on standard benchmark pose estimation datasets such as LINEMOD [17], Occluded-LINEMOD [18], and YCB-Video [13], but they require either intensive manual annotation or a motion capture system. To the authors’ knowledge, a readily applicable method for collecting underwater training data with the corresponding accurate 6D poses is not available. In this work, we focus on estimating the 6D pose of an Aqua2 vehicle [19] (shown in Fig. 1). The observer is either another Aqua2 robot or an underwater handheld camera. The proposed method utilizes the Unreal Engine 4 (UE4) [20] with a 3D model of the Aqua2 robot swimming, projected over underwater images to generate training images with known poses for the pose estimation network. Dissimilarity in images arising from intrinsic factors such as distortion differences from different cameras, external factors such as color-loss, poor visibility quality, or the surroundings, hampers the performance of

classical deep learning-based 6D pose estimation methods in the underwater domain. CycleGAN [21] was employed to transform UE4 rendered images to image sets used for training with varying in appearance, similar to real-world underwater images.

Using a modified version of YOLOv3 [22] to detect an object bounding box, the proposed network produces robust 6D pose estimates by combining multiple local predictions of 2D keypoints that are projections of 3D corners of the object. Only grid cells inside the detected bounding box contribute to the selection of 2D keypoints along with a confidence score. Using the predictions with confidence, the most dependable 2D keypoint candidates for each 3D keypoint are selected to yield a set of 2D-to-3D correspondences. These selected 2D keypoints are used in the RANSAC-based PnP [23] algorithm to obtain a robust 6D pose estimate.

The proposed framework has been tested in different environments – pool, ocean – and different platforms, including an Aqua2 robot and GoPro cameras, demonstrating its robustness. The main contributions are as follows:

- A 6D pose prediction network that predicts object bounding boxes and eight keypoints in image coordinates. These 2D keypoints are then used in 2D-to-3D correspondence to estimate 6D pose.
- We demonstrate the effective use of rendered image augmentation² in 6D pose prediction, eliminating the need for ground truth labeling in real images. Utilizing image augmentation from the rendered to the underwater environment, the pose prediction network becomes invariant to color-loss, texture-smoothing, and other domain-specific challenges.
- We publish a dataset of the Aqua2 robot captured in the ocean and swimming pool to further research in the underwater domain³.

The next section reviews related works. Section III introduces the proposed method, including the synthetic data generated for training, and the pose estimation method. Section IV presents first the ground truth data acquisition used exclusively for testing, then quantitative results from different datasets together with a comparison with other methods are discussed. Finally, we conclude the paper with future work in section V.

II. RELATED WORK

In this paper, we focus on 6D pose estimation using RGB images without access to a depth map. RGB-D based methods [24]–[26] are not applicable underwater given the attenuation of infrared light at a very short distance. The classical approach for 6D object pose estimation involves extracting local features from the input image, matching them with features from a 3D model to establish 2D-to-3D correspondences from which a 6D pose can be obtained through the PnP algorithm [7], [8], [10], [27]. Previous

work studied local feature descriptors invariant to changes in scale, rotation, illumination, and viewpoints [28], [29]. Even though these feature-based techniques can handle occlusions and scene clutter, they require sufficient texture to compute local features. To deal with poorly-textured objects, some efforts focused on learned feature descriptors using machine learning techniques [30], [31].

In recent years, pose estimation research has been dominated by frameworks utilizing deep neural network. These methods can be broadly classified into two categories: either regressing directly to 6D pose estimates [11], [13], or predicting 2D projections of 3D keypoints on an image and then obtaining pose via PnP algorithm [14], [15]. Xiang *et al.* [13] estimate the object center in the image with the distance of the center used for estimating the translation and the predicted quaternions for object rotations. Peng *et al.* [32] used a pixel-wise voting network to regress pixel-wise unit vectors pointing to the keypoints and used these vectors to vote for keypoint locations using RANSAC. Recent works [33]–[35] researched on post-processing to refine the initial pose estimates from the first step. Li *et al.* [36] disentangled the pose to predict rotation and translation separately from two different branches to increase accuracy.

Recent approaches [37], [38] for pose estimation focus on local patches belonging to the object rather than producing a single global prediction. The work of Hu *et al.* [16] is closest to our approach in terms of using local image patches which also learns a semantic segmentation mask to select multiple keypoint locations from local patches belonging to an object and providing those inputs to the PnP algorithm. Regarding pose estimation using synthetic datasets, Rozantsev *et al.* [39] used a two-stream network trained on a synthetic and real dataset, and introduced loss functions that prevent corresponding weights of two streams from being too different from each other. Rad *et al.* [40] proposed a method that learns a feature mapping from real to synthetic datasets, and during inference transfers the features of real images to synthetic and infers pose using synthetic features. Some work has been done using a deep learning framework for Aqua2 vehicle detection that enabled visual servoing [4]. Koreitem *et al.* [41] used rendered images for pose estimation based visual tracking of Aqua2, and our approach outperforms their approach in terms of 6D pose estimation accuracy.

Our work employs CycleGAN [21], a type of Generative Adversarial Network (GAN), to generate a synthetic dataset for training. GANs, introduced by Goodfellow *et al.* [42], are used to generate images through adversarial training where a generator attempts to produce realistic images to fool a discriminator which tries to distinguish if the image is real or generated. CycleGAN [21] is used for unpaired image-to-image translation even in absence of corresponding real-generated image pair. The main idea of CycleGAN is that if an image is translated from one domain to another and translated back, the resulting image should resemble the original image.

²Traditionally Generative Adversarial Networks (GAN) [21] use the term image translation for this operation, however, the term translation can be confusing for a robotics application.

³<https://afrl.cse.sc.edu/afrl/resources/datasets/>

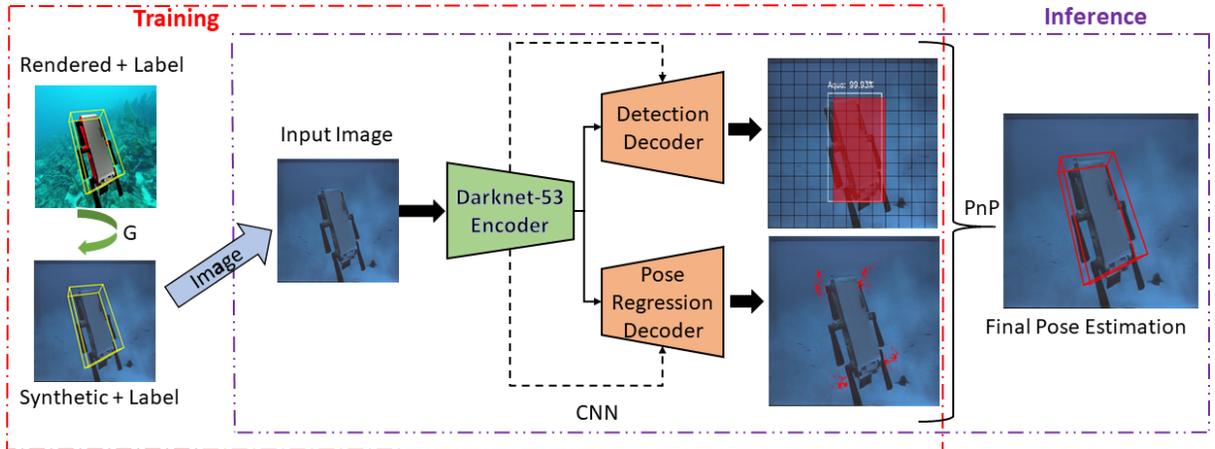


Fig. 2: In training (outlined in red), the rendered images are translated to the synthetic images resembling Aqua2 swimming in a pool or a ocean environment. The synthetic images are then fed to a common encoder, which is connected to two decoder streams: Detection Decoder (object detection) and Pose Regression Decoder (6D pose regression). Only in inference (outlined in purple), are the predicted 2D keypoint projections of 8 corners of the 3D Aqua2 model processed and utilized to obtain a 6D pose using the RANSAC-based PnP algorithm.

III. THE PROPOSED SYSTEM

Figure 2 shows an overview of the proposed system. In the training process, UE4 renders a 3D model of Aqua2 with known 6D poses projected on top of underwater ocean images. The feature space between real underwater and rendered images is aligned by transferring the rendered images to target domains (swimming pool and ocean) using CycleGAN [21], an image-to-image translation network.

The next stage consists of a Convolutional Neural Network (CNN) that predicts the 2D projections of the 8 corners of the object’s (Aqua2) 3D model, similar to [14], [15] and an object detection bounding box. Even though [14], [15] divide an image into grid cells, they use global estimates of 2D keypoints for the object with the highest confidence value. In our approach, each grid cell inside the bounding box predicts the 2D projections of keypoints along with their confidences focusing on local regions belonging to the object. These predictions of all cells are then combined based on their corresponding confidence scores using RANSAC-based PnP during 6D pose estimation.

A. Domain Adaptation

We employ CycleGAN [21] for unpaired image-to-image translation by learning functions to map the UE4 domain R to the target domain T . We use generators G and F to transfer domains: $G : R \rightarrow T$ and $F : T \rightarrow R$. Discriminator, D_R , is designed to distinguish between rendered images in R , and augmented fake images $F(T)$. Discriminator, D_T , aims to separate target images in T and augmented fake images $G(R)$. To improve image-to-image translation in CycleGAN, cycle consistency is maintained by ensuring the reconstructed images $F(G(R)) \approx R$ in addition to the adversarial loss. To calculate adversarial loss, G tries to generate $G(R)$, which is so similar to T that can fool the discriminator D_T . The

loss for G and D_T is:

$$L_G(G, D_T, R, T) = E_{t \sim p_{\text{data}}(t)}[\log D_T(t)] + E_{r \sim p_{\text{data}}(r)}[\log(1 - D_T(G(r)))] \quad (1)$$

where $t \sim p_{\text{data}}(t)$ and $r \sim p_{\text{data}}(r)$ denotes the data distribution in T and R respectively, and E is the loss function, which is L1-norm in our approach. Similarly we derive $L_R(F, D_R, T, R)$ following Eq. 1. The cycle consistency loss L_{cyc} is defined as:

$$L_{\text{cyc}}(G, F) = E_{r \sim p_{\text{data}}(r)}[\|F(G(r)) - r\|_1] + E_{t \sim p_{\text{data}}(t)}[\|G(F(t)) - t\|_1] \quad (2)$$

In our proposed method, there are two target domains: swimming pool, T_{sp} and an open-water ocean environment, T_m . Therefore, we train two instances of CycleGAN (two generators), $G_1 : R \rightarrow T_{sp}$ and $G_2 : R \rightarrow T_m$. Fig. 3 shows the CycleGAN training overview along with synthetic data generation.

B. 6D Pose Estimation

The proposed network consists of an encoder, Darknet-53 [22], and two decoders: Detection Decoder and Pose Regression Decoder. The detection decoder detects objects with bounding boxes, and the pose regression decoder regresses to 2D corner keypoints of the 3D object model. The decoders predict the output as a 3D tensor with a spatial resolution of $S \times S$ and a dimension of D_{det} and D_{reg} , respectively. The spatial resolution controls the size of an image patch that can effectively vote for object detection and for the 2D keypoint locations. The feature vectors are predicted at three different spatial resolutions. The decoder stream detects features with multiple scales via upsampling and concatenation with a depth of final layer, D_{det} . The pose regression stream also has a similar architecture, but the final depth layer is maintained to be D_{reg} . Predicting in

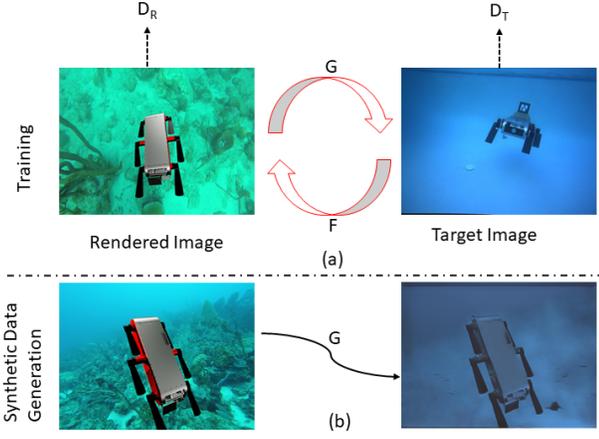


Fig. 3: (a) CycleGAN learning process is shown. CycleGAN learns two mapping functions; $G : R \rightarrow T$ and $F : T \rightarrow R$ with two discriminators, D_R and D_T . (b) Only using generator G , we perform image-to-image translation of rendered images R to target images T .

multiple spatial resolutions with upsampling helps to obtain semantic information at multiple scales using fine-grained features from early on in the network.

Object Detection Stream: The object detection stream is similar to the detection stream of YOLOv3 [22] which predicts object bounding box. For each grid cell at offset (c_x, c_y) from the top left corner of the image, the network predicts 4 coordinates for each bounding box p_x, p_y, p_w, p_h . Following [22], we use 9 anchor boxes obtained by k-means clustering on COCO dataset [43] of size $(10 \times 13), (16 \times 30), (33 \times 23), (30 \times 61), (62 \times 45), (59 \times 119), (116 \times 90), (156 \times 198), (373 \times 326)$ divided among three scales. The width and height are predicted as the fraction of the anchor box priors a_w, a_h and the actual bounding box values are obtained as

$$\begin{aligned} b_x &= \sigma(p_x) + c_x \\ b_y &= \sigma(p_y) + c_y \\ b_w &= a_w e^{p_w} \\ b_h &= a_h e^{p_h} \end{aligned} \quad (3)$$

where σ represents the sigmoid function. The sum of square of error between the ground truth t_* and coordinate prediction \hat{t}_* is used as the loss function. The ground truth values t_* can be obtained by inverting equation Eq. (3). The object detection stream also predicts the objectness score of each bounding box by calculating its intersection over union with anchor boxes and class prediction scores using independent logistic classifiers as in [22]. The total object detection loss L_{det} is the sum of coordinate prediction loss, objectness score loss, and class prediction loss. The total object detection loss was introduced by Redmon *et al.* [44] to which we refer for a complete description.

Pose Regression Stream: The pose regression stream predicts the location of the 2D projections of the predefined 3D keypoints associated with the 3D object model of Aqua2. We

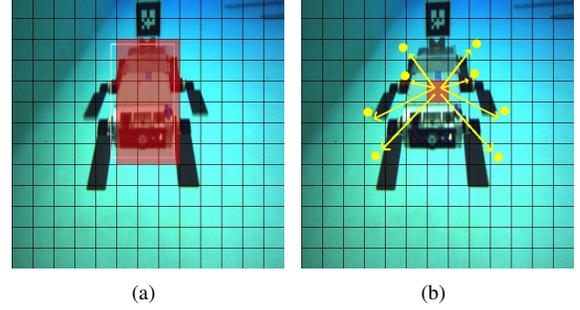


Fig. 4: (a) The object detection stream predicts the bounding box and assigns each cell inside the box to the Aqua2 object. (b) The regression stream predicts the location of 8 bounding box corners as 2D keypoints from each grid cell.

use 8 corner points of model bounding boxes as keypoints. The pose regression stream predicts a 3D tensor with size $S \times S \times D_{\text{reg}}$. We predict the (x, y) spatial locations for the 8 keypoint projections along with their confidence values, $D_{\text{reg}} = 3 \times 8$.

We do not predict the 2D coordinates of the 2D keypoints directly. Rather, we predict the offset of each keypoint from the corresponding grid cell as in Fig. 4(b) in the following way: Let c be the position of grid cell from top left image corner. For the i^{th} keypoint, we predict the offset $f_i(c)$ from grid cell, so that the actual location in image coordinates becomes $c + f_i(c)$, which should be close to the ground truth 2D locations g_i . The residual is calculated as

$$\Delta_i(c) = c + f_i(c) - g_i \quad (4)$$

and we define offset loss function, L_{off} , for spatial residual:

$$L_{\text{off}} = \sum_{c \in B} \sum_{i=1}^8 \|\Delta_i(c)\|_1 \quad (5)$$

where B consists of grid cells that fall inside the object bounding box and $\|\cdot\|_1$ represents L1-norm loss function, which is less susceptible to outliers than L2 loss. Only using grid cells falling inside the object bounding box for 2D keypoint predictions focuses on image regions that truly belong to the object.

Apart from the 2D keypoint locations, the pose regression stream also calculates the confidence value $v_i(c)$ for each predicted point, which is obtained through the sigmoid function on the network output. The confidence value should be representative of the distance between the predicted keypoint and ground truth values. A sharp exponential function of the 2D euclidean distance between prediction and ground truth is used as confidence. The confidence loss is calculated as

$$L_{\text{conf}} = \sum_{c \in B} \sum_{i=1}^8 \|v_i(c) - \exp(-\alpha \|\Delta_i(c)\|_2)\|_1 \quad (6)$$

where $\|\cdot\|_2$ denotes euclidean distance or L2 loss and parameter α defines the sharpness of the exponential function. The pose regression loss of Eq. (8) takes up the form

$$L_{\text{reg}} = \lambda_{\text{off}} L_{\text{off}} + \lambda_{\text{conf}} L_{\text{conf}} \quad (7)$$

For numerical stability, we down-weight the confidence loss for cells that do not contain objects by setting λ_{conf} to 0.1, as suggested in [44]. For the cells that include the object, λ_{conf} is set to 5.0 and λ_{off} to 1. Therefore, the total loss of the network is:

$$L = L_{\text{det}} + L_{\text{reg}} \quad (8)$$

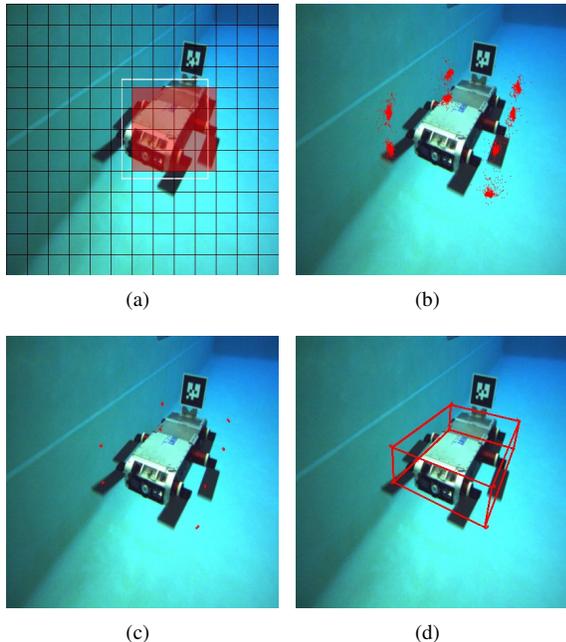


Fig. 5: Inference strategy for combining pose candidates. (a) Grid cells inside the detection box belonging to Aqua2 object overlaid on the image. (b) Each grid predicts 2D locations for corresponding 3D keypoints shown as red dots. (c) For each keypoints, 12 best candidates are selected based on the confidence scores. (d) Using $12 \times 8 = 96$ 2D-to-3D correspondence pairs and running RANSAC-based PnP algorithm yield accurate pose estimate as shown by the overlaid bounding box.

C. Pose Refinement

During inference, the object detection stream of our network predicts the coordinate locations of the bounding boxes with their confidences and the class probabilities for each grid cell. Then, the class-specific confidence score is estimated for the object by multiplying the class probability and confidence score. To select the best bounding box, we use non-max suppression [45] with an IOU threshold of 0.4 and a class-specific confidence score threshold of 0.3.

Simultaneously, the pose regression stream produces the projected 2D locations of the object’s 3D bounding box, along with their confidence scores for each grid cell, as shown in Fig. 5-b. The 2D keypoint predictions for grid cells that fall outside of the bounding box (Fig. 5-a) from the object detection stream are filtered out. In an ideal case, the remaining 2D keypoints should cluster around the object center. 2D keypoints that do not belong to a cluster are removed using a pixel distance threshold of 0.3 times image width.

The keypoints with confidence scores less than 0.5 are also filtered out. To balance the trade-off between computation time and accuracy, we empirically found that using the 12 most confident 2D predictions for each 3D keypoint (Fig. 5-c) produces an acceptable pose estimate after RANSAC-based PnP [23]. Hence, we employ RANSAC-based PnP [23] on $12 \times 8 = 96$ 2D-to-3D correspondence pairs between the image keypoints and the object’s 3D model to obtain a robust pose estimate, as shown in Fig. 5-d.

D. Implementation Details

To create the synthetic dataset, we train the CycleGAN following the training procedure of [21]. We let the training continue until it generated acceptable reconstruction. Once CycleGAN can reasonably reconstruct for the target domain, we use the model weights of that epoch to translate all rendered images to synthetic images. Then, the synthetic images are scaled to 416×416 resolution maintaining the aspect ratio by padding zeros for training. During inference, no augmentation is required, and the real images are directly fed to the network.

The CNN is trained for 125 epochs on the synthetic dataset, and the first 3 epochs are part of a warmup phase, where the learning rate gradually increases from 0 to $1e-4$. We utilized the SGD optimizer with a momentum of 0.9 and a piecewise decay to decrease the learning rate to $3e-5$ and $1e-5$ at 60 and 100 epochs, respectively. To avoid overfitting, minibatches of size 8 were produced by applying data augmentation techniques, including randomly changing hue, saturation, and exposure of the image up to a factor of 1.5. In addition, images were randomly scaled, and affine transformed by up to 25% of the original image size.

IV. EXPERIMENTS

This section describes first the datasets used, and then results of the inference with the real Aqua2 robot swimming in both a pool and the open ocean at Barbados.

A. Datasets Description

Training - Rendered/Synthetic Dataset: contains images obtained by rendering an Aqua2 robot swimming with flipper motion using UE4 and overlaying the resulted 3D model over random underwater images. Rather than just overlaying the 3D model of Aqua2, we simulate the flipper motion to generate images with the flippers in various realistic positions. This flipper motion makes the neural network independent of the flipper position. The synthetic dataset is obtained by the image-to-image translation network based on CycleGAN described in III-A to create photo-realistic images. The rendered dataset contains 37K images with random depth between 0.75 m and 3.0 m and orientations ranging from -50 to 50 degrees for roll, -70 to 70 degrees for pitch, and -90 to 90 degrees for yaw.

Testing - Pool Dataset: To generate our pool dataset, we deployed two robots: one robot observing the other with a vision-based 2D fiducial marker (AR tag⁴) mounted

⁴http://wiki.ros.org/ar_track_alvar

	Translation Error	Orientation Error	REP-10px Accuracy	ADD-0.1d Accuracy	FPS
Tekin <i>et al.</i> [15]	0.278m	18.87°	9.33%	23.39%	54
PVNet [32]	0.486m	24.55°	23.22%	43.09%	37
DeepURL	0.068m	6.77°	25.22%	57.16%	40

TABLE I: Translation and Orientation errors (the lower the better) along with REP-10px, ADD-0.1d accuracy (the higher the better) and runtime comparison for the pool dataset

on the top used to estimate ground truth during two pool trials in indoor and outdoor pools, as shown in Fig. 6(a-d). Approximately, 11K images were collected with estimated localization, provided from the pose detection of the AR tag and the relative transformation of the mounted tag to the real robot. The dataset contains images with a distance between two Aquas ranging between 0.5m to 3.5m.

Testing - Barbados 2017 Dataset: The Barbados 2017 Dataset consists of 188 real images collected during underwater field trials off the west coast of Barbados used in [41], see Fig. 6(i-l). The images are captured from an Aqua2 robot’s onboard camera. 6D pose of the robot in each of these images is obtained using a custom-built annotator, which allows the user to mark keypoints on the robot assigned from the CAD model. The annotator then iteratively fits a wireframe to the robot using its known dimensions.

Testing - Barbados GoPro Dataset: We collected images underwater in Barbados of an Aqua2 robot swimming over coral reefs using a GoPro camera, which differ significantly from the images collected using another Aqua2 in terms of hue, image size, and aspect ratio (see Fig. 6(e-h)). Given that ground truth is unavailable for these images, this dataset was only used to evaluate the proposed method qualitatively.

B. Evaluation Metrics

To evaluate the pose estimation capability of the proposed system, we calculated the mean translation error as the Euclidean distance between the predicted and the ground truth translation. Let $(Rot, trans)$ and (\hat{Rot}, \hat{trans}) be the ground truth and predicted rotation matrices and translation, respectively. For individual angle errors in terms of yaw, roll, and pitch, we decomposed the rotation matrices Rot and \hat{Rot} into Euler angles and calculated their absolute difference. The total orientation error is represented as Eq. (9), where tr represents the trace of the matrix and the orientation error is in the range of $[0, \pi]$.

$$\phi(Rot, \hat{Rot}) = \arccos \frac{tr(Rot^T \hat{Rot}) - 1}{2} \quad (9)$$

To evaluate the pose accuracy, we use standard metrics - namely- 2D reprojection error [46] and the average 3D distance of the model vertices, referred to as ADD metric [13], [14]. In the case of reprojection error, we consider the pose estimate as correct if the average distance between 2D projections of 3D model points obtained using predicted and ground-truth poses is below a 10 pixels threshold, referred to as REP-10px. Generally, a 5-pixel threshold is employed, but we consider a threshold of 10 pixels to account for

uncertainties in ground truth due to the AR tag-based pose estimation. The ADD metric takes pose estimate as correct if the mean distance between the coordinates of the 3D model vertices transformed by estimated and ground truth pose fall below 10% of the model diameter, referred to as ADD-0.1d. We also report the inference time of the algorithm in terms of frames per second (FPS) on an RTX 2080 GPU.

C. Experimental Results

Evaluation on the Pool Dataset: We compare our method with the state-of-the-art method of Tekin *et al.* [15] and PVNet [32] trained on a synthetic dataset and tested on real pool dataset. Translation and rotation errors along with REP-10px and ADD-0.1d accuracy for the pool dataset are presented in Table I, as well as the runtime comparisons on Nvidia RTX 2080. DeepURL outperforms both the method of Tekin *et al.* [15] and PVNet [32] in terms of rotation and translation errors along with REP-10px and ADD-0.1d accuracy. Moreover, the runtime performance is realtime, outperforming PVNet [32] and only slightly slower than that of Tekin *et al.* [15]. The improved performance comes from two factors: 1) the use of a better detection pipeline and 2) bounding box based keypoint sampling introduced in this paper. Whereas Tekin *et al.* [11] only used the keypoints with the highest confidence, our bounding box based keypoint sampling allows the selection of more appropriate keypoints using RANSAC-based PnP. PVNet [32], compared to DeepURL, performed slightly inferior on REP-10px and ADD-0.1d metrics and produced significantly higher translation and orientation errors.

Figure 7 shows the translation and orientation error statistics of an Aqua2 robot in the pool dataset. It is evident that the proposed method performs well across all distances from camera (0.5m-3.5m). Interestingly, at very close distance, the method experiences higher orientation error due to the 2D keypoints of Aqua2 not being precisely selected by the pose regression decoder.

Evaluation on the Barbados 2017 Dataset: We report the performance of our system on Real Barbados 2017 dataset in terms of translation and rotation errors as shown in Table II. DeepURL performs significantly better on translation error and orientation error compared to the method of Koreitem *et al.* [41].

Impact of Domain Adaptation: To understand the efficacy of using CycleGAN based domain adaptation, we trained DeepURL only on rendered images. Even though the network performed well on the validation set consisting of rendered images only, without training on synthetic dataset it was not able to generalize to real-world pool images.

	Translation Error	Orientation Error	Roll Error	Pitch Error	Yaw Error
Koreitem <i>et al.</i> [41]	0.72m	17.59°	11.87°	4.59°	12.11°
DeepURL	0.31m	11.98°	9.64°	3.30°	5.43°

TABLE II: Translation and Rotation errors for the Barbados 2017 dataset [41]

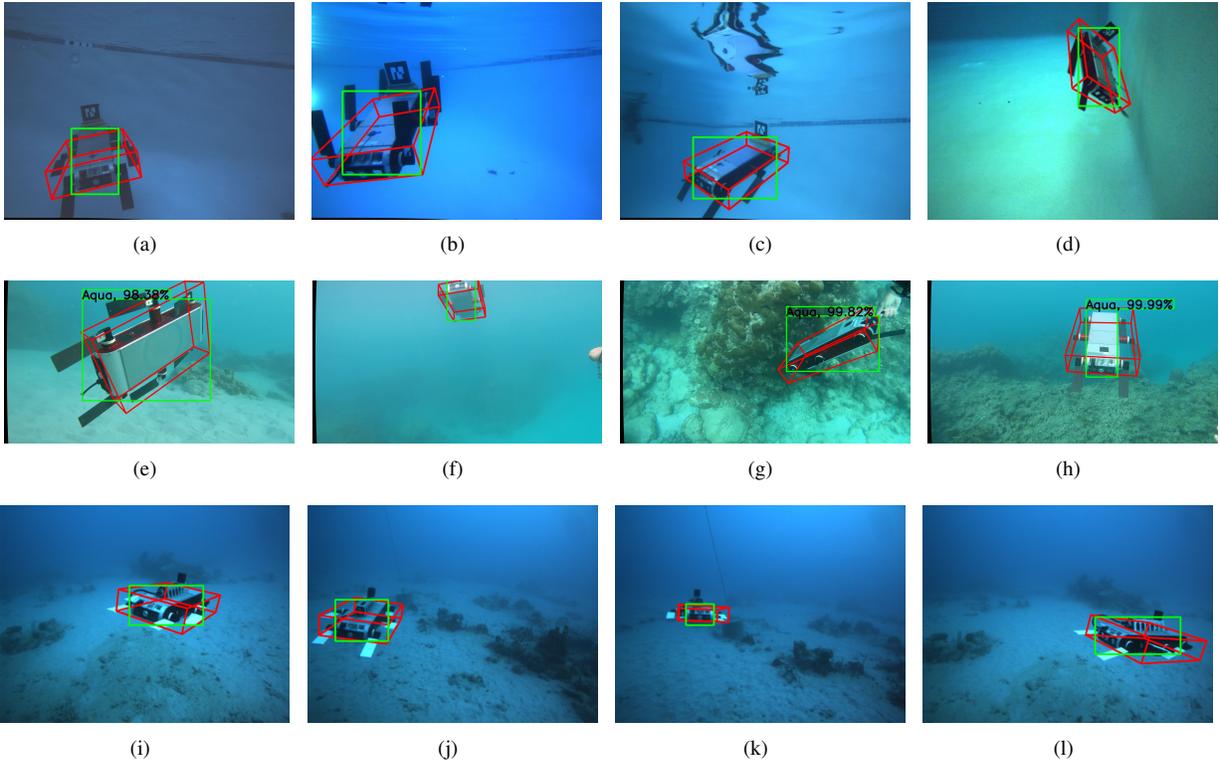


Fig. 6: Sample detections from the different datasets. Green square is the 2D detection box, while the red wireframe is the projection of the 3D bounding box of the robot. Top row: Pool dataset, observed Aqua2 vehicle carries a AR tag to generate ground truth estimates; observing robot is another Aqua2. Middle row: GoPro footage during deployments in Barbados in January 2020, observed robot has no additional components, the observing camera is a GoPro 7 camera. Bottom row: Barbados 2017 dataset [41], observed robot is equipped with a Ultra-Short Baseline (USBL) modem, the observing robot is an Aqua2 vehicle.

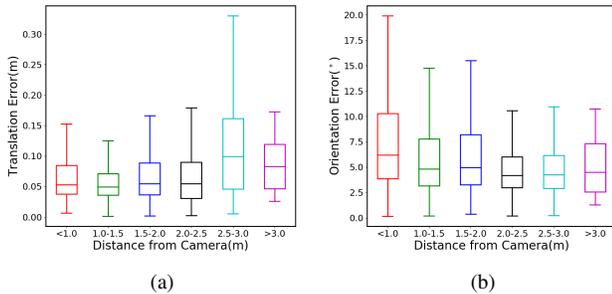


Fig. 7: Boxplot summarizing the error statistic of (a) translation and (b) orientation in Pool Dataset with respect to variable distance from camera.

The intuition is that the real-world underwater images differ significantly from rendered images in terms of texture, color and appearance. Thus, creating image sets with different appearance and texture helps extensively in the training process by reducing over-fitting and increasing generalization.

Failing Scenarios: Predictions of pose estimate might be wrong either when the detection stream fails to predict the object detection box, therefore, there is not enough points for the RANSAC-based PnP algorithm (at least six 2D-to-3D correspondences are required), or PnP did not

converge. These detection failure scenarios are inherent in YOLOv3 architecture. The system may also fail for position or orientation not introduced in the training scenarios, such as translation beyond 3.5m or orientation beyond the rendered range described in Section IV-A.

V. CONCLUSION

In this work, we presented a system for 6D pose estimation of an autonomous underwater vehicle for relative localization underwater. The system learns to predict the 6D pose without the need for any real ground truth, which enables pose estimation in an environment where ground truth is difficult to acquire. We also present a detection bounding box based keypoint sampling strategy that is more robust to related work [15], [32] which leads to a better estimate of the pose of the observed robot, up to an order of magnitude in some cases; see Table I.

Currently, the proposed network is being ported to an Intel Neural Compute Stick 2 (Intel NCS2)⁵ and an NVidia Jetson TX2 Module⁶ in order to deploy on an Aqua2 or a BlueROV2 vehicle. The above two platforms were selected based on their performance [47] and compatibility with the

⁵<https://software.intel.com/en-us/neural-compute-stick>

⁶<https://developer.nvidia.com/embedded/jetson-tx2>

proposed vehicles. Furthermore, the DeepURL framework will be integrated with the proprioceptive sensors of each robot (IMU and depth) and either the USBL positioning of the observer or the Visual-Inertial estimator [48] to recover the pose of both robots in a global frame of reference.

REFERENCES

- [1] I. M. Rekleitis, G. Dudek, and E. E. Milios, "On multiagent exploration," in *Vision Interface*, 1998, pp. 455–461.
- [2] A. Quattrini Li, A. Coskun, S. M. Doherty, S. Ghasemlou, A. S. Jagtap, M. Modasshir, S. Rahman, A. Singh, M. Xanthidis, J. M. O’Kane, and I. Rekleitis, "Experimental comparison of open source vision based state estimation algorithms," in *Proc. ISER*, 2016, pp. 775–786.
- [3] B. Joshi, S. Rahman, M. Kalaitzakis, B. Cain, J. Johnson, M. Xanthidis, N. Karapetyan, A. Hernandez, A. Quattrini Li, N. Vitzilaios, and I. Rekleitis, "Experimental Comparison of Open Source Visual-Inertial-Based State Estimation Algorithms in the Underwater Domain," in *Proc. IROS*, 2019, pp. 7221–7227.
- [4] F. Shkurti, W. Chang, P. Henderson, M. Islam, J. Higuera, J. Li, T. Manderson, A. Xu, G. Dudek, and J. Sattar, "Underwater multi-robot convoying using visual tracking by detection," in *Proc. IROS*, 2017, pp. 4189–4196.
- [5] T. Manderson, J. Li, N. Dudek, D. Meger, and G. Dudek, "Robotic Coral Reef Health Assessment Using Automated Image Analysis," *J. Field Robot.*, vol. 34, no. 1, pp. 170–187, 2017.
- [6] T. Manderson, J. C. Gamboa, S. Wapnick, J.-F. Tremblay, D. Meger, and G. Dudek, "Vision-based goal-conditioned policies for underwater navigation in the presence of obstacles," in *Proc. RSS*, July 2020.
- [7] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proc. ICCV*, vol. 2, 1999, pp. 1150–1157.
- [8] F. Rothganger, S. Lazebnik, C. Schmid, and J. Ponce, "3D object modeling and recognition using local affine-invariant image descriptors and multi-view spatial constraints," *Int. J. Comput. Vision*, 2006.
- [9] A. Collet, M. Martinez, and S. S. Srinivasa, "The MOPED framework: Object recognition and pose estimation for manipulation," *Int. J. Robot. Res.*, vol. 30, no. 10, pp. 1284–1306, 2011.
- [10] D. Wagner, G. Reitmayr, A. Mulloni, T. Drummond, and D. Schmalstieg, "Pose tracking from natural features on mobile phones," in *IEEE/ACM Int. Symp. on Mix. and Aug. Reality*, 2008, pp. 125–134.
- [11] W. Kehl, F. Manhardt, F. Tombari, S. Ilic, and N. Navab, "SSD-6D: Making RGB-based 3D detection and 6D pose estimation great again," in *Proc. ICCV*, 2017.
- [12] M. Sundermeyer, Z.-C. Marton, M. Durner, M. Brucker, and R. Triebel, "Implicit 3D orientation learning for 6d object detection from rgb images," in *ECCV*, 2018.
- [13] Y. Xiang, T. Schmidt, V. Narayanan, and D. Fox, "PoseCNN: A convolutional neural network for 6D object pose estimation in cluttered scenes," *Proc. RSS*, 2018.
- [14] M. Rad and V. Lepetit, "BB8: a scalable, accurate, robust to partial occlusion method for predicting the 3d poses of challenging objects without using depth," in *Proc. ICCV*, 2017.
- [15] B. Tekin, S. N. Sinha, and P. Fua, "Real-time seamless single shot 6d object pose prediction," in *Proc. CVPR*, 2018.
- [16] Y. Hu, J. Hugonot, P. Fua, and M. Salzmann, "Segmentation-driven 6d object pose estimation," in *Proc. CVPR*, 2019.
- [17] S. Hinterstoisser, V. Lepetit, S. Ilic, S. Holzer, G. Bradski, K. Konolige, and N. Navab, "Model based training, detection and pose estimation of texture-less 3D objects in heavily cluttered scenes," in *ACCV*, 2013, pp. 548–562.
- [18] A. Krull, E. Brachmann, F. Michel, M. Ying Yang, S. Gumhold, and C. Rother, "Learning analysis-by-synthesis for 6d pose estimation in rgb-d images," in *Proc. ICCV*, 2015.
- [19] G. Dudek, M. Jenkin, C. Prahacs, A. Hogue, J. Sattar, P. Giguere, A. German, H. Liu, S. Saunderson, A. Ripsman, S. Simhon, L. A. Torres-Mendez, E. Milios, P. Zhang, and I. Rekleitis, "A visually guided swimming robot," in *Proc. IROS*, 2005.
- [20] T. Manderson, I. Karp, and G. Dudek, "Aqua underwater simulator," in *Proc. IROS*, 2018.
- [21] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *Proc. ICCV*, 2017.
- [22] J. Redmon and A. Farhadi, "YOLOv3: An incremental improvement," 2018.
- [23] V. Lepetit, F. Moreno-Noguer, and P. Fua, "EPnP: An accurate O(n) solution to the PnP problem," *Int. Journal of Comp. Vision*, 2009.
- [24] E. Brachmann, A. Krull, F. Michel, S. Gumhold, J. Shotton, and C. Rother, *Learning 6D Object Pose Estimation Using 3D Object Coordinates*. Springer, 2014, vol. 8690.
- [25] C. Choi and H. I. Christensen, "RGB-D object pose estimation in unstructured environments," *Robotics and Aut. Systems*, vol. 75, pp. 595 – 613, 2016.
- [26] J. Sock, S. H. Kasaei, L. S. Lopes, and T.-K. Kim, "Multi-view 6d object pose estimation and camera motion planning using rgbd images," *IEEE Int. Conf. on Comp. Vision Workshops (ICCVW)*, 2017.
- [27] G. Pavlakos, X. Zhou, A. Chan, K. G. Derpanis, and K. Daniilidis, "6-DoF object pose from semantic keypoints," in *Proc. ICRA*, 2017, pp. 2011–2018.
- [28] D. Lowe, "Distinctive image features from scale-invariant keypoints," *IJCV*, 11 2004.
- [29] T. Trzcinski, M. Christoudias, V. Lepetit, and P. Fua, "Learning image descriptors with the boosting-trick," in *Proc. NeurIPS*, 2012.
- [30] P. Wohlhart and V. Lepetit, "Learning Descriptors for Object Recognition and 3D Pose Estimation," in *Proc. CVPR*, 2015.
- [31] A. Doumanoglou, V. Balntas, R. Kouskouridis, and T.-K. Kim, "Siamese regression networks with efficient mid-level feature extraction for 3D object pose estimation," 2016.
- [32] S. Peng, Y. Liu, Q. Huang, X. Zhou, and H. Bao, "PVNet: Pixel-wise Voting Network for 6DoF Pose Estimation," in *Proc. CVPR*, 2019.
- [33] S. Zakharov, I. Shugurov, and S. Ilic, "DPOD: 6D Pose Object Detector and Refiner," in *Proc. ICCV*, 2019.
- [34] K. Gupta, L. Petersson, and R. Hartley, "CullNet: Calibrated and Pose Aware Confidence Scores for Object Pose Estimation," in *ICCV Workshops*, 2019.
- [35] Y. Li, G. Wang, X. Ji, Y. Xiang, and D. Fox, "DeepIM: Deep Iterative Matching for 6D pose estimation," in *Proc. ECCV*, 2018.
- [36] Z. Li, G. Wang, and X. Ji, "CDPN: Coordinates-Based Disentangled Pose Network for Real-Time RGB-Based 6-DoF Object Pose Estimation," in *Proc. ICCV*, 2019.
- [37] O. Hosseini Jafari, S. K. Mustikovela, K. Pertsch, E. Brachmann, and C. Rother, "iPose: Instance-aware 6D pose estimation of partly occluded objects," in *Computer Vision – ACCV*, 2018, pp. 477–492.
- [38] M. Oberweger, M. Rad, and V. Lepetit, "Making deep heatmaps robust to partial occlusions for 3D object pose estimation," *CoRR*, vol. abs/1804.03959, 2018.
- [39] A. Rozantsev, M. Salzmann, and P. Fua, "Beyond sharing weights for deep domain adaptation," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 41, pp. 801–814, 2019.
- [40] M. Rad, M. Oberweger, and V. Lepetit, "Feature mapping for learning fast and accurate 3D pose inference from synthetic images," in *Proc. CVPR*, 2018.
- [41] K. Koreitem, J. Li, I. Karp, T. Manderson, F. Shkurti, and G. Dudek, "Synthetically trained 3d visual tracker of underwater vehicles," in *MTS/IEEE OCEANS*, Charleston, SC, USA, Oct. 2018.
- [42] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Proc. NeurIPS*, 2014, pp. 2672–2680.
- [43] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft COCO: Common objects in context," in *Proc. ECCV*, 2014, pp. 740–755.
- [44] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proc. CVPR*, 2016.
- [45] R. Rothe, M. Guillaumin, and L. V. Gool, "Non-maximum suppression for object detection by passing messages between windows," in *ACCV*, 2014.
- [46] E. Brachmann, F. Michel, A. Krull, M. Ying Yang, S. Gumhold, and c. Rother, "Uncertainty-driven 6D pose estimation of objects and scenes from a single RGB image," in *Proc. CVPR*, 2016.
- [47] M. Modasshir, A. Quattrini Li, and I. Rekleitis, "Deep neural networks: a comparison on different computing platforms," in *Proc. CRV*, 2018, pp. 383–389.
- [48] S. Rahman, A. Quattrini Li, and I. Rekleitis, "SVIn2: An Underwater SLAM System using Sonar, Visual, Inertial, and Depth Sensor," in *Proc. IROS*, 2019, pp. 1861–1868.