

**SENSOR BASED TARGET TRACKING WITH APPLICATION TO AUTONOMOUS
DOCKING AND SELF-RECONFIGURABILITY**

Shubhdildeep S. Sohal

Robotics and Mechatronics Lab
Mechanical Engineering Dept.
Virginia Tech
Blacksburg, VA, USA
shubh94@vt.edu

Pinhas Ben-Tzvi *

Robotics and Mechatronics Lab
Mechanical Engineering Dept.
Virginia Tech
Blacksburg, VA, USA
bentzvi@vt.edu

ABSTRACT

This paper presents a target detection technique, which combines a supervised learning model with sensor data to eliminate false positives for a given input image frame. Such a technique aids with selective docking procedures where multiple robots are present in the environment. Hence the sensor data provides additional information for this decision making process. Sensor accuracy plays a crucial role when the motion of the robot is defined by the use of data recorded by its sensors. The uncertainties in the sensory data can cause misalignments due to poor calibration of the sensor, which can result in poor positioning of the robot relative to its target. Such misalignments can play a significant role where certain accuracy is desired. Therefore, it is necessary to minimize such misalignments to achieve certainty for the robot interaction with its target. The work proposed in this paper allows achieving such accuracy using a vision-based approach by eliminating all false occurrences leading to selective interactions with the target. The proposed methodology is validated using a self-reconfigurable mobile robot capable of hybrid Wheeled-Tracked mobility, as an application towards autonomous docking of mobile robotic modules.

Keywords: Autonomous Docking, Image processing, Target detection and tracking, Self-reconfigurable robots.

1 INTRODUCTION

The autonomy of robots has been on the rise with an increasing demand for versatile, intelligent, and adaptable robotic structures [1-5]. Modern-day robots are equipped with multiple sensors to collect data from the surrounding for an intelligent decision-making process. An example of such robots has been shown in Fig. 1(a) and 1(b), where the symmetrically invertible robots are equipped with onboard sensors to help the

robot navigate and interact in an unknown environment. These autonomously operating robots can be assigned to move from point A to point B with a mere interaction of onboard sensors with the robots' surroundings. In the last decade, the use of Neural Networks has been substantial [6-8] and has emerged as one of the key focus areas in the field of computer vision, particularly as an application for object detection, classification, tracking, etc. (herein objects have been referred to as targets). Using such an approach, these robots can be programmed to quickly make intelligent decisions to avoid any obstacles whether they are static or dynamic to reach their destination. These networks can be trained to solve the problems of feature detection and image classification. The amount of information available with the image is of great importance while designing such networks as it can increase the network complexity and

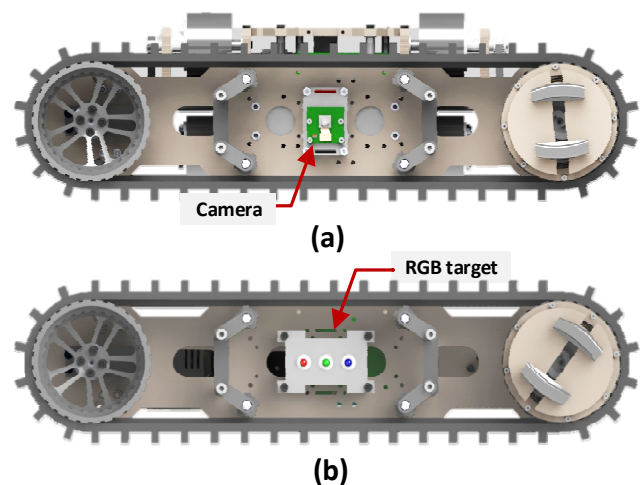


Figure 1. Symmetrically invertible modular robots (a) Locomotion robot, (b) Manipulator robot

computation cost. Recently, proposed networks have shown a great prospect in application to the field of robotics due to their reliability with vision-based sensors [6-8].

To begin with the proposed approach, a motivating application involving multi-robot assembly has been shown in Fig. 2 [9-11]. A Self Transformable Robotic Module (STORM) has been developed [12], which consists of two main parts, i.e., Hybrid-Wheeled locomotion robot and a manipulator robot as shown in Fig. 2. In the example shown in Fig. 2, the combined formation consists of a manipulator robot and two locomotion robots which consist of coupling mechanisms called GHEFT (Genderless High strength Efficient Fail-safe and high misalignment Tolerant) [13]. This selection has been highlighted based on the review of the coupling mechanism in [14] and due to its misalignment tolerance capacity along X-Y-Z-roll-pitch-yaw axes of (6, 25, 11) mm and (45, 11, 11) degrees respectively. The 2-DOF genderless docking mechanism allows for independent clamping translation and the rotation of the mechanism. These symmetrically invertible robots can operate under flip-over conditions. The locomotion robot is equipped with a Hybrid-Wheeled assembly capable of exhibiting bi-directional mobility along the longitudinal and lateral direction. As part of a future application, it is proposed that the two locomotion robots will approach a manipulator robot from either side using path planning techniques as shown in Fig. 2 (e.g. using Pozyx sensor due to their better accuracy compared to other position based sensors such as GPS, Bluetooth, Wi-Fi) [15-17]. However, these sensors fail to provide an absolute positioning of the target which fails to dock under autonomous locomotion. To align the robots within the misalignment tolerance range of the docking GHEFT mechanism, a vision-based approach is used. The inaccuracies related to these sensors can be used to estimate the Region of Interest (*RoI*) which further will be combined with the detection technique to extract the location of the target in the image plane.

This approach serves as detection speed and accuracy improvement over the Hybrid-Target Tracking technique proposed in [18]. The error in the positioning sensor accounts for the use of such detection techniques due to its lower

complexity and faster tracking with better accuracy. Since this technique lacks the ability to handle scalability, pose, and variable lighting of the target, a new approach has been presented in this research, which is discussed in detail in later sections. The proposed approach can also be combined with the autonomous landing of drones using GPS based beacons, or RTK-GPS based positioning along with the use of vision-based markers.

The paper starts with the introduction of the proposed approach followed by a discussion of the experimental setup and the results. Lastly, a brief conclusion summarizing the approach in this paper and directions for future work has also been presented.

2 PROPOSED APPROACH

The recent shift in the use of Neural Networks to classify different objects is a significant improvement over the conventional sliding window techniques such as Template Matching. The use of such networks allows achieving higher accuracy and precision with the detection even under changing scalability and lighting conditions. A classic example would be to detect the handwritten digits using LeNet-5 [7]. It is an 8-layer structure consisting of an input layer, two convolutional layers, two sub-sampling layers, three fully connected layers, and an output layer. Here the number of output neurons varies from 0 to 9.

The recent development in this research is the inclusion of YOLO (You Only Look Once) [8] object classifier, capable of detecting 80 classes with a mean average precision (mAP) of 57.9%. These networks highly depend on the amount of training being used to generate the model. To achieve high accuracy with detections, such a network is equipped with a large number of convolutional layers, which increases the complexity and computation of the network. The implementation YOLO is not feasible on small portable micro-computers such as Raspberry Pi due to their limited computation power. However, a small version of YOLO called YOLO-tiny presents much better compatibility on these devices. Such a network takes approximately 24-25 secs to process an image, which can be further reduced to 2-3 secs

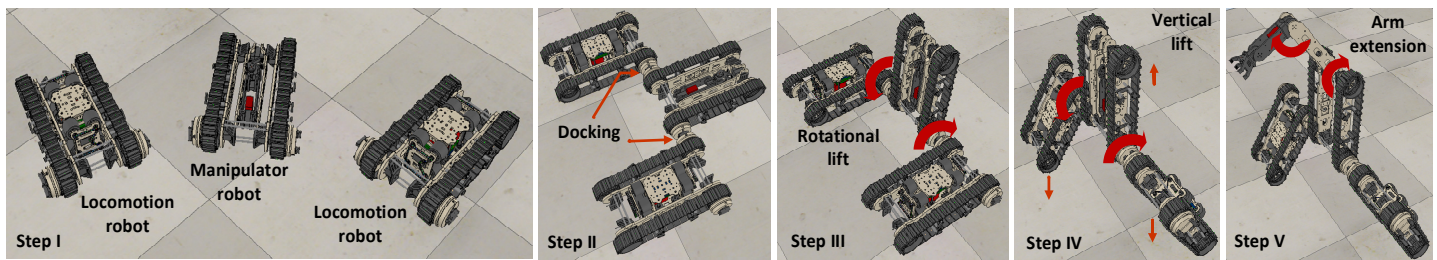


Figure 2.3 Robot docking shown in step-wise simulation for humanoid configuration performed in V-REP with STORM Locomotion and Manipulator module.

using NNPACK. NNPACK is an acceleration package for Neural Network computations for multi-core CPUs. Such an approach becomes redundant when we already have some estimation of the target position in the image plane. To highlight these issues, the paper presents the following two contributions:

1. Integration of the sensor data with the detection methodology to minimize the layer density of the CNN (Convolutional Neural Networks), which in turn speeds up the detection.
2. Using sensor data to present the selective detection for autonomous docking in mobile robots as an application for the proposed methodology. Such a method will help to eliminate the extra True-Positive (TP) and unnecessary False-Positive (FP) detections. Thus, it helps to achieve a higher True Positive Rate ($TPR = TP / (TP + FN)$) and Positive Predictive Value (precision or $PPV = TP / (TP + FP)$). Here, the term FN represents the False-negative value. The detection behavior has been shown in Fig. 3, where multiple detections including both TP and FP detection formed on an image plane, relating to the detection of the target robots as seen from the source robot.

The methodology proposed in this paper aims at using sensor data to reduce the computational load of CNN. A brief layout of the process flow-chart is shown in Fig. 4. The common available CNN models classify each object in terms of a class, such that a given frame could have n possible detections of the same class. Such methodology is limited to a search over the whole frame since there is no sensorial information available with the detected objects, e.g. person, car, etc. The proposed approach is validated by using a target attached to the side frame of the robot and using another robot as the initializer for the detection. The camera is attached to a Hybrid-Wheeled mobile robot (source), whereas the colored target is attached to the side frame of the manipulator robot (target). The proposed version of the algorithm takes into account the uneven behavior of the terrain with an added external tilt value along the Roll and Pitch axes. The final performance of the methodology is estimated by the time it

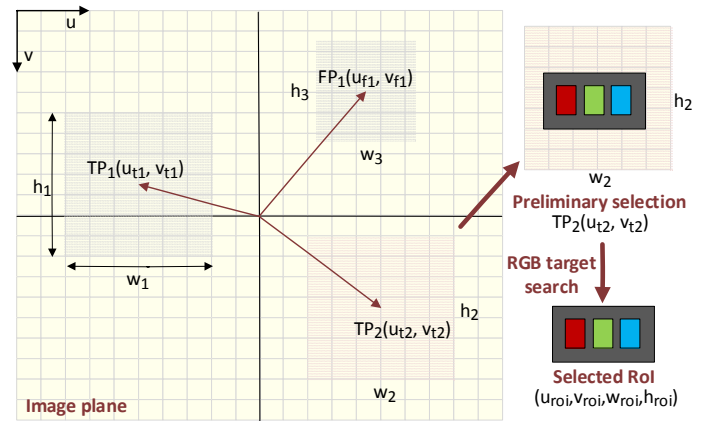


Figure 3. Presumed occurrences of the targets on the image plane showing the True Positive (TP), False Positive (FP) and desired detections (red area).

takes to process and correctly identify the target in the image.

The proposed methodology is an extension of the previous work mentioned in [18]. The previous version of the target tracking technique suffered from challenges such as scalability and changing illumination in the background. Moreover, it was further assumed that the target is always within the image plane as viewed by the camera. Furthermore, only one target was placed in the image frame as a proof of concept of the algorithm to avoid any false positive detections. Since phase 1 of the algorithm provides the initial estimate of the detected target, the algorithm needs to provide a fast and reliable scalability invariant detection, in the presence of occlusion and light variability. Therefore, this paper proposes the use of a supervised learning model to counter the failures with phase 1 detection of the previous algorithm [18]. There are several state-of-the-art supervised learning models available in the literature; however, they have network layers to detect all the possible trained object appearances at high accuracy. To do that, an image is used as an input to the detection model, where convolutional, pooling, activation, and other functions are applied to the image. In cases where only one object is present in an image, applying the above-mentioned operations on the whole image will be computationally expensive and highly complex. The inclusion of the sensor data to determine the

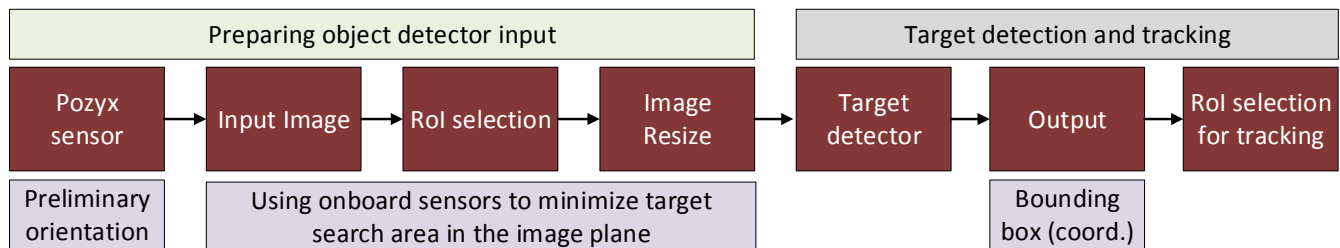


Figure 4. Detection process before initializing the target tracking.

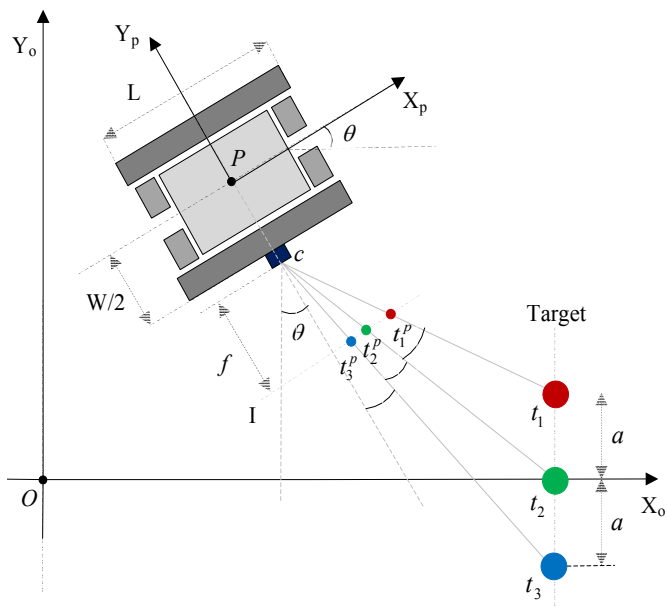


Figure 5. Target features placed at a certain distance from the robot are projected over the image plane.

relative pose and approximate position of the target helps in eliminating the false positive by minimizing the input region to the detection model. Furthermore, such an approach helps in minimizing the layer density/complexity of the network.

The proposed approach takes advantage of the Pozyx sensory data available with the target robot(s) to differentiate the FP from TP detections. The positioning sensor can provide a 2D and 3D positioning along with an onboard IMU sensor based on the reference recorded from the 4 anchors placed within the operating range. As shown in Fig. 3 if three targets are detected using a conventional CNN based approach, then the approach used in these applications eliminates the FP_1 detection leaving the robot with only two TP detections (TP_1 and TP_2). The detected two TP targets are further differentiated using the sensor data from each target robot. Since each detected TP target will be identified by its sensory data, the docking becomes much easier in terms of target selectivity. This selectivity is shown in Fig. 3, marked by a red region, selecting TP_2 out of TP_1 and TP_2 . Each TP detection of the marked region comprises an RGB LED colored marker, such that the camera attached to the side frame of the source robot is used to track the motion relative to the target robot in the Image Plane. An example of such a scenario has been shown in Fig. 5, where RGB colored markers represent the marker attached to the side frame of the target robot and top view of the source robot representing the detection using the onboard camera.

The selected region is used as an input to find the position of the target (u_{target}, v_{target}) relative to the image plane ($u_{roi}, v_{roi}, w_{roi}, h_{roi}$) based on the trained network model.

$$\begin{aligned} u_{target} &= (u_2 - \frac{w_2}{2} + u_{roi} + \frac{w_{roi}}{2}) \\ v_{target} &= (v_2 - \frac{h_2}{2} + v_{roi} + \frac{h_{roi}}{2}) \end{aligned} \quad (1)$$

Since each occurrence of the target can vary in terms of scale and orientation, it is required to accommodate such variations using an additional approximation from the positioning sensors. To achieve this flexibility, the detected bounding box parameters (u_p, v_p, w, h), are then used as the input for phase 2 of the algorithm defined by parallel tracking with the optical flow ($u_p + \Delta u_p, v_p + \Delta v_p, t + \Delta t$) and box segmentation for the LED (LED_1, LED_2, LED_3) color tracking. This optical flow process is shown in Fig. 6, where the input represents the bounding box coordinates given by the target detector and output represents the individual tracking of the colored markers using the color segmentation approach. The t value defines the update in the pixel coordinates of the tracked point over time. This flow estimate of the tracked coordinate is summarized as follows,

$$I(u_p, v_p, t) = I(u_p + \Delta u_p, v_p + \Delta v_p, t + \Delta t) \quad (2)$$

These target LEDs are represented by a red-blue-green color combination ($t_1 - t_2 - t_3$) as shown in Fig. 5. The projection of these colored targets on the image plane is given by $t_1^p - t_2^p - t_3^p$ and in terms of pixel coordinates as (u_1, v_1) , (u_2, v_2) , and (u_3, v_3) . The use of optical flow technique in parallel to color tracking provides the following advantages:

1. It delivers consistent tracking of the target robot, even though the color tracking may fail at a certain time-step providing consistency with the performance at the same time.

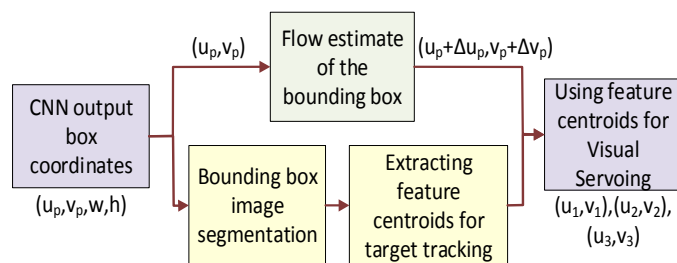


Figure 6. Target detector output for the phase 2 of the Hybrid Target Tracking (HTT) technique.

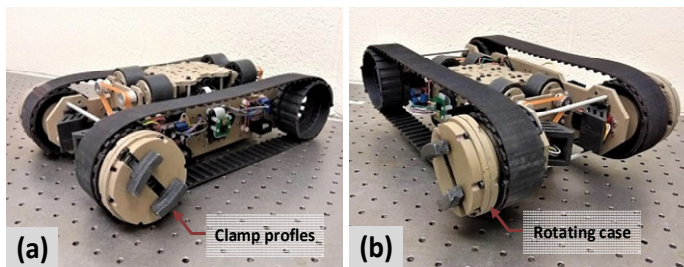


Figure 7. STORM robot interfaced with GHEFT docking mechanism, showing (a) track-actuated mode, (b) Wheeled mode for multi-directional locomotion

2. The use of limited RoI helps to improve the performance speed and reduces the search (or processing) area for color tracking compared to a search over the entire image.
3. The optical flow estimate of the RoI also accommodates the scalability of the target, in cases where the source robot approaches (or moves away from) the target robot.

A basic convolutional neural network consists of an input layer, an output layer, and multiple hidden layers (convolution and fully connected layers). The proposed model uses the raw image data based on the extracted RoI and generates the output in terms of the location of the desired target in that image. Thus, the output is defined in terms of the following parameters: width (w) and height (h) of the bounding box, pixel coordinates of the location, and the confidence value of the prediction. The aim is to minimize the number of filters by reducing the feature map size.

3 EXPERIMENTAL SETUP

The experimental setup to implement the algorithm on the mobile robot is discussed in this section.

The training of the model was done using a host PC equipped with i5 Processor, 8GB RAM, and NVIDIA GTX 1060 GPU. To validate the proposed algorithm, it was run on a Hybrid-Wheeled Mobile Robot equipped with an onboard Raspberry Pi 3 Model B computer (1GB RAM) and an IMU sensor to record the roll (α), pitch (β), and yaw (γ) orientation of the robot. An Arducam 5MP camera was attached to the side frame of the source robot and an RGB colored marker was attached to the side frame of the target robot. A 400x400 image was captured to initialize the detection process.

The bi-directional locomotion mechanism of the Hybrid-Wheeled mobile robot, as shown in Fig. 7, allows for the switch in mobility based on the locomotion requirement on different terrains. The mechanism, which controls the bi-directional mobility of the robot is shown in Fig. 8(a). Fig. 8(b) shows the cut-section view of the robot, representing the Hybrid-Wheeled assembly attached to the track-actuated side frame of the robot. This bi-directional mobility allows the robot to move in the longitudinal direction using tracked locomotion mode and lateral direction using the wheeled locomotion mode. Additionally, the robots are equipped with a Pozyx sensor to get the positioning data (x, y, z) of the target robot with sub-meter accuracy, and an onboard IMU sensor to get the roll-pitch-yaw values ($\alpha_s, \beta_s, \gamma_s$) and ($\alpha_t, \beta_t, \gamma_t$), of the source and target robot respectively. Here the subscript s and t represent the source and target robot respectively. The sensor data is transferred from the target robot to the source robot using an onboard Wi-Fi module on the Raspberry Pi. Apart from the sensors, the robot is capable of mobility (along with longitudinal and lateral directions) with an actuation of a prismatic joint. The robot is interfaced with a 2-DOF docking mechanism, which can tolerate mechanical misalignments along X-Y-Z and about the

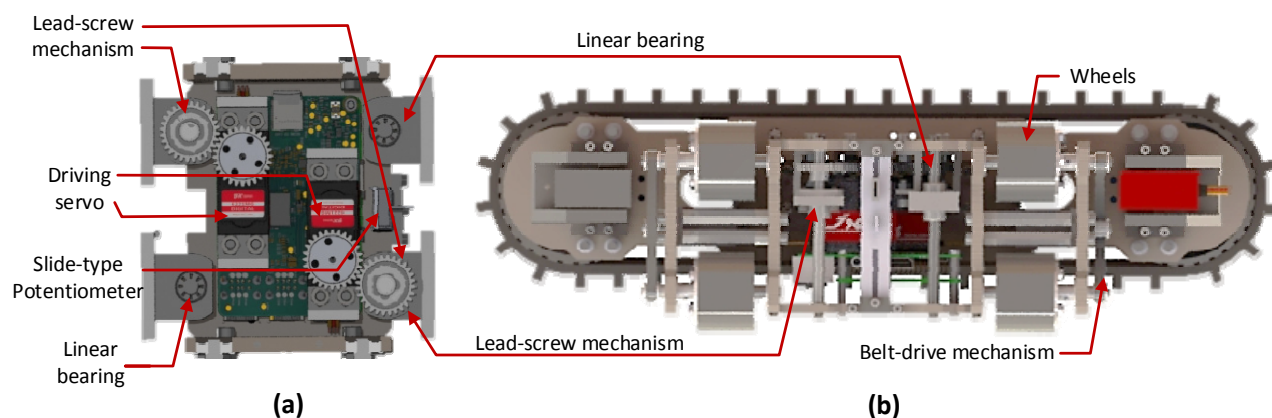


Figure 8. Cut-section of Vertical Translational Unit which enables multi-directional mobility of the robot (a) dual drive mechanism using driving motors on either side for a synchronized motion, (b) prismatic joint between the Hybrid-Wheeled assembly and the side frame of the driving mechanism.

roll-pitch-yaw axes.

$$(\Delta x, \Delta y, \Delta z) = \min\{(x_s, y_s, z_s), (x_t, y_t, z_t)\} \quad (3)$$

Here the (x, y, z) represents the positioning data while the s and the t subscript represent the source and the target robot respectively. The source robot was placed in close vicinity to the target robot considering the error along X-Y-Z axes based on the data given by the Pozyx sensors. In such a case, the minimization of the error using the IMU sensor data from the source and target robot can be given as

$$(\Delta\alpha, \Delta\beta, \Delta\gamma) = \min\{(\alpha_s, \beta_s, \gamma_s), (\alpha_t, \beta_t, \gamma_t)\} \quad (4)$$

Here (α, β, γ) represents the roll-pitch-yaw angles recorded using onboard IMU data while s and t subscript represents the source and the target robot respectively. These error angle values are subject to fall within certain limit values defined as

$$(\Delta\alpha, \Delta\beta, \Delta\gamma) = \begin{cases} \alpha_{l1} < \Delta\alpha < \alpha_{l2} \\ \beta_{l1} < \Delta\beta < \beta_{l2} \\ \gamma_{l1} < \Delta\gamma < \gamma_{l2} \end{cases} \quad (5)$$

These limit values $(\alpha_{l1}, \alpha_{l2}, \beta_{l1}, \beta_{l2}, \gamma_{l1}, \gamma_{l2})$ are defined relative to the target robot. In case there are multiple robots in the environment, these values will adhere relatively to the target robot with which the source robot will interact. The initial orientation alignment for facing the docking mechanism of the source and the selected target robot can be done using Eq. (3) and Eq. (4). The limit values are used to accommodate the inaccuracy with both the positioning and the orientation of the robot. The initial RoI selection based on the sensor data is done by projecting the target pixel coordinates using the extrinsic

parameters given by the Pozyx sensor $(\Delta x, \Delta y, \Delta z)$ and IMU data $(\Delta\alpha, \Delta\beta, \Delta\gamma)$, combined with the intrinsic parameters of the camera (f, ρ, u_o, v_o) . These intrinsic parameters correspond to the focal length, pixel size, and the image center of the camera. The width (w_{roi}) and the height (h_{roi}) of the RoI are calculated using triangular similarity based on ΔZ , f and the actual dimensions $(w_t$ and $h_t)$ of the target. The real-world positioning helps in providing an approximate location estimate of the target in the image. Moreover, it helps to minimize the search area to be used as an input for the detection model. The detection of the RoI is followed by the optical flow technique for consistent tracking of the target.

Apart from these listed error minimizations, an additional experimentation result has been shown in the following section to demonstrate the autonomous locomotion of the source robot relative to the target robot.

4 EXPERIMENTAL RESULTS AND ANALYSIS

The section presents a discussion related to the autonomous motion of the self-reconfigurable shown in Fig. 9 and Fig. 10. Figure 9 represents the use of the proposed methodology on the locomotion robot, which is further utilized to track the target features in the Image Plane. Figure 9(a) shows the robot equipped with the Pozyx sensor to get the X-Y-Z position of the target relative to the source robot. The trajectory of the robot has been shown in Fig. 9(b) in red. The robot is set to move from its initial position to the final position which is the center of the Image Plane. The motion of each colored feature on the Image plane is represented by red-blue-green curves along the u and v axis, respectively, using Image-Based Visual Servoing (IBVS) [19]. The motion of the features helps to demonstrate the convergence of the error between the initial pixel coordinates and the desired location (herein center

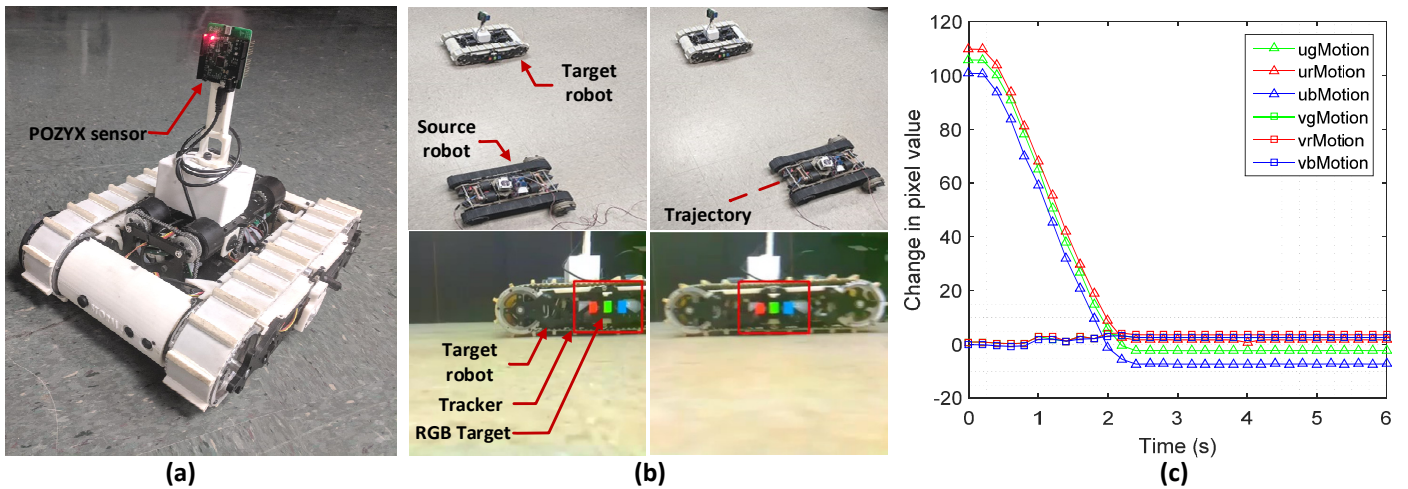


Figure 9. (a) STORM locomotion robot equipped with Pozyx sensor, (b) Motion of the target features using Visual Servoing, (c) Curves representing u and v pixel motion for the three LED features.

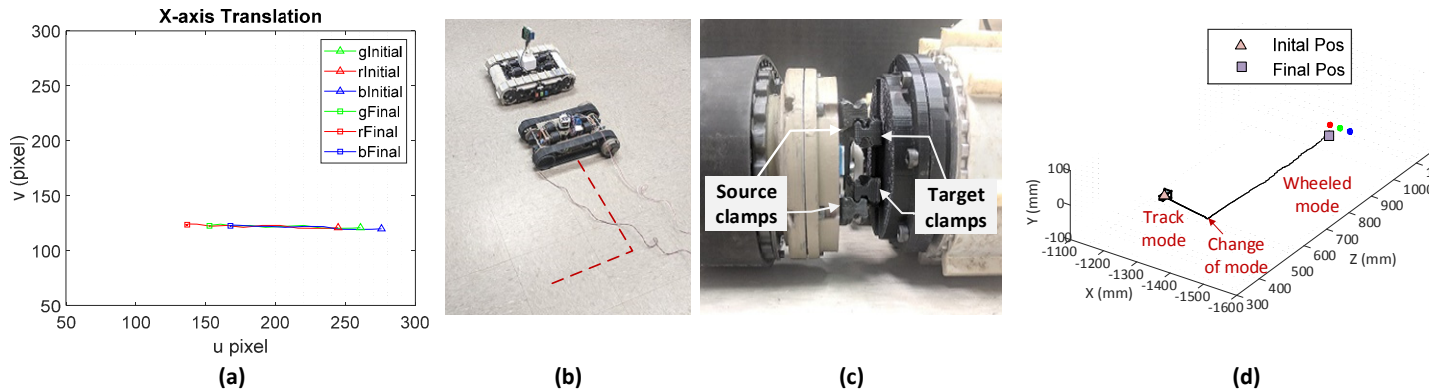


Fig. 10. (a) Feature error plots for the motion of the source robot along X, Y, Z and Yaw direction relative to the target robot, (b) Motion trajectory of the source robot represented in red color, (c) Clamping between the robots for self-reconfigurability, (d) Real world motion trajectory tracking of the robot.

of the Image Plane). There is a minor shift in the pixel coordinates of the features along the v axis which can be related to a minor vibration of the robot. The corresponding motion of the target features in the image plane, as viewed from the camera attached to the side frame of the robot is shown in Fig. 10 (a). This translation motion along the X-axis will be followed by a sequence of motions along the Y-axis and the Z-axis respectively (shown in Fig. 10(b)) to couple the docking mechanism of the source robot to that of the target robot as shown in Fig. 10(c). The final motion trajectory of the source robot relative to the target robot is shown in Fig. 10(d), which is generated using the onboard sensor data.

Moreover, it should be noted that the current experimentation only included basic preliminary testing on a vinyl floor. The impact of vibrations can be high depending on the change in roughness or unevenness of the surface, e.g. gravel, grass, etc. Since this work is a part of ongoing research, the testing related to the target tracking on different surfaces [16] will be demonstrated as a part of future research.

5 CONCLUSION AND FUTURE WORK

This work presented a new approach in object tracking, which combines the sensor data with a supervised learning model to improve the detection accuracy of the targets and to provide flexibility for selective target tracking. The approach utilized sensor data to restrict the search area in the input image, thereby reducing the complexity and the difficulty of training the network. The availability of sensor data with the target serves as a preliminary estimate in this technique. The proposed methodology is implemented on a self-reconfigurable mobile robot for autonomous docking using IBVS in an indoor environment. However, the proposed technique will be further used to test the autonomous docking of the robotic modules in an outdoor environment using other sensors, which can provide better accuracy such as RTK-GPS.

Future work involves outdoor testing of the robotic modules on uneven terrain to validate the proposed approach and the self-reconfigurability of mobile robots. The development of more robots will also help to generate a multi-robot coordination network, which will help to illustrate the docking selectivity of the method. Furthermore, other networks can be analyzed to draw a performance comparison between the proposed approach and the existing methods. The effect of dimensionality reduction of the features in an image will also be analyzed in the future.

REFERENCES

- [1] Zhong, M., Guo, W., Li, M., Xu, J., "Tanbot: A Mobile Self-Reconfigurable Robot Enhanced with Embedded Positioning Module", 2008 IEEE Workshop on Advanced Robotics and It's Social Impacts, Taipei, 2008, pp. 1-5.
- [2] Daudelin, J., Jing, G., Tosun, T., Yim, M., Gazit, H. K., Campbell, M., "An integrated system for perception-driven autonomy with modular robots", IEEE Science Robotics, 2018. doi: 10.1126/scirobotics.aat4983
- [3] Ben-Tzvi, P., Goldenberg, A. A., and Zu, J. W., "Articulated hybrid mobile robot mechanism with composed mobility and manipulation and onboard wireless sensor/actuator control interfaces," *Mechatronics Journal*, vol. 20, no. 6, pp. 627-639, Sept. 2010.
- [4] Murata, S., Kakomura, K., and Kurokawa, H., "Docking Experiments of a modular robot by visual feedback," *International Conference on Intelligent Robots and Systems*, Beijing, 2006, pp. 625-630.
- [5] Moubarak, P., P. Ben-Tzvi, "Adaptive Manipulation of a Hybrid Mechanism Mobile Robot", *Proc. of the 2011 IEEE International Symposium on Robotic and Sensors Environments (ROSE 2011)*, Canada, pp. 113-118, Sept. 17-18, 2011.
- [6] Ren, X. D., Guo, H. N., He, G. C., Xu, X., Di, C., Li, S. H., "Convolutional Neural Network Based on Principal Component Analysis Initialization for Image

- Classification”, IEEE First International Conference on Data Science in Cyberspace (DSC), Changsha, pp. 329-334, 2016. doi: 10.1109/DSC.2016.18.
- [7] Lecun, Y., Bottou, L., Bengio, Y., Haffner, P., “Gradient-based learning applied to document recognition”, Proceedings of the IEEE, vol. 86, no. 11, pp. 2278-2324, 1998.
- [8] Redmon, J., Divvala, S., Girshick, R., Farhadi, A., “You Only Look Once Unified, Real-time Object Detection”, IEEE Conference on Computer Vision and Pattern Recognition. 2016.
- [9] Moubarak, P., Ben-Tzvi, P., "A Tristate Rigid Reversible and Non-Back-Drivable Active Docking Mechanism for Modular Robotics," IEEE/ASME Transactions on Mechatronics, Vol. 19, Issue 3, pp. 840-851, June 2014.
- [10] Moubarak, P. M., and Ben-Tzvi, P.,” On the Dual-Rod Slider Rocker Mechanism and its applications to Tristate Rigid Active Docking,” Journal of Mechanisms and Robotics, Vol. 5, Issue 1, pp. 011010:1-10, Feb. 2013.
- [11] Moubarak, P. M., Alvarez, E. J., and Ben-Tzvi, P.,” Reconfiguring a Modular Robot into a Humanoid Formation: A Multi-Body Dynamic Perspective On Motion Scheduling for Modules and Their Assemblies”, Proc. of 2013 IEEE Int. Conf. on Autom. Sci. and Eng., Madison, Wisconsin, Aug. 17-21, 2013.
- [12] Kumar, P., Saab, W., Ben-Tzvi, P., “Design of a Multi-Directional Hybrid-Locomotion Modular Robot with Feedforward Stability Control”, Proceedings of the 2017 ASME IDETC/CIE, 41st Mechanisms & Robotics Conference, Ohio, Aug. 6-9, 2017.
- [13] Saab, W., Ben-Tzvi, P.,” A Genderless Coupling Mechanism with 6-DOF Misalignment Capability for Modular Self-Reconfigurable Robots”, Journal of Mechanisms and Robotics, Transactions of the ASME, Vol. 8, Issue 6, pp. 061014:1-9, Dec. 2016.
- [14] Saab, W., Racioppo, P., Ben-Tzvi, P., “A review if coupling mechanism design for modular reconfigurable robots”, Robotica Journal, October 2018. doi:10.1017/S0263574718001157.
- [15] Sebastian, B., Ben-Tzvi, P., “Physics-Based Path Planning for Autonomous Tracked Vehicles in Challenging Terrain”, Journal of Intelligent and Robotic Systems, 2018. doi: 10.1007/s10846-018-0851-3.
- [16] Sebastian, B., Ben-Tzvi, P., “Active disturbance rejection control for handling slip in tracked vehicle locomotion”, Journal of Mechanisms and Robotics, Dec. 2018. doi:10.1115/1.4042347.
- [17] Pozyx sensors for positioning and motion information. Available. [Online]. <https://www.pozyx.io/>
- [18] Sohal, S. S., Saab, W., Ben-Tzvi, P., “Improved Alignment Estimation for Autonomous Docking in Mobile Robots”, Proc. of the 2018 ASME IDETC/CIE, 42nd Mechanisms and Robotics Conference, Quebec City, Canada, Aug. 26-29, 2018.
- [19] Corke, P., “MATLAB toolboxes: robotics and vision for students and teachers”, IEEE Robotics and Automation Magazine, vol. 14, Issue 4, pp. 16-17, 2007.