# A Closed-Loop Multi-perspective Visual Servoing Approach with Reinforcement Learning

Lei Zhang[1,2]†, Jiacheng Pei[2,3]†, Kaixin Bai[1,2], Zhaopeng Chen[2,1]*, Jianwei Zhang[1]

*Abstract*—Traditional visual servoing methods suffer from serving between scenes from multiple perspectives, which humans can complete with visual signals alone. In this paper, we investigated how multi-perspective visual servoing could be solved under robot-specific constraints, including self-collision, singularity problems. We presented a novel learning-based multi-perspective visual servoing framework, which iteratively estimates robot actions from latent space representations of visual states using reinforcement learning. Furthermore, our approaches were trained and validated in a Gazebo simulation environment with connection to OpenAI/Gym. Through simulation experiments, we showed that our method can successfully learn an optimal control policy given initial images from different perspectives, and it outperformed the Direct Visual Servoing algorithm with mean success rate of 97.0%.

Fig. 1. The sythetic robot agent is controlled with proposed closed-loop visual servoing approach to target pose with desired visual state.

## I. INTRODUCTION

Humans can guide their behavior using semantic information from visual images captured at different angles. Similarly, visual servoing enables robots to adjust their motion based on visual feedback. Visual servoing is influenced by several factors, including feature matching and robot trajectory planning. However, traditional visual servoing methods often assume that the initial observation pose is similar to the target's observation pose, limiting their applicability in multi perspective visual servoing scenarios.

Multi-perspective visual servoing in complex industrial scenarios poses greater challenges due to difficulties in acquiring accurate object models and the inherent complexity of scenes. Occlusions of objects also sharpen the difficulties of pose estimation, making 6D pose-based visual servoing inadequate for such situations. In contrast, traditional visual servoing methods are expert in achieving high-speed, low-error robot guidance within local regions based on image features or specific features. However, precise servoing to the target position from arbitrary poses remains challenging. Hand-crafted features are often used in traditional visual servoing methods, which leads to traditional visual servoing being applicable only to small convergence domains. Significant disparities between the image features of the initial state and the target state result in increased challenges
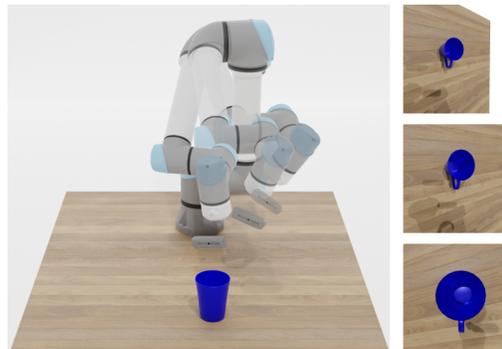
in feature matching of visual servoing. Simultaneously, the constraints of the robot need to be taken into account during visual servoing. Ultimately, the difficulties in multi-view visual servoing persist due to perceptual limitations and the constraints imposed by the robot.

Recently learning-based visual servoing has also been studied. Deep learning-based approaches have shown enhanced performance in dealing with complicated scenarios [1], occluded environments, and varying lighting conditions [2], [3]. Reinforcement learning were also utilized in improving generalization of visual servoing [4]–[9]. However, most of the current learning-based visual servoing approaches only consider 2D perspectives and do not address the problem of robotic operations under robotic constraints. Training robust robot control strategies under robot-specific constraints are also the core problem of deep learning-based robot manipulation methods. In this work, we consider the problem of multi-perspective visual servoing using reinforcement learning (RL) methods, as shown in Fig. 1. Our method utilizes autoencoder network to extract latent space representations of current and desired camera sensor data. We propose a closed-loop robot control policy network to estimate robot action from the latent space representations. Our core contributions are:

- A closed-loop multi-perspective visual servoing framework utilizing RL to servo robot agent from latent space representations of visual states.
- Improve training efficiency of RL-based robotic policy using learning from demonstration method and Hindsight Experience Replay (HER) [10].
- A potential-based reward function with consideration of the task- and robot-specific constraints.

†The first two authors contribute equally to this paper.

*Corresponding author.

[1]TAMS (Technical Aspects of Multimodal Systems), Department of Informatics, University of Hamburg, [2]Agile Robots AG, [3]RWTH Aachen University

## II. RELATED WORK AND BACKGROUND

### A. Visual Servoing

Image-based visual servoing (IBVS) [11], [12] uses the extracted 2D features as states to achieve robot control, whereas pose-based visual servoing (PBVS) [13], [14] estimates the poses of camera in the Cartesian space and guides the robot by minimizing pose error. IBVS suffers from convergence and stability problems due to ill-conditioned image Jacobian matrix and singularities [15] and the problem of finding optimal path in the Cartesian space [16]. PBVS is limited by the image quality and calibration errors. The traditional methods require tracking a set of handcrafted features, such as points, lines or patterns. Direct or photometric visual servoing (DVS) treats the luminance of image pixels as the visual features and computes the interaction matrix [17]–[21]. However, DVS suffers from small convergence domain. It may be failed when features are not visible or in challenging lighting condition [2].

Recent advances in machine learning open up new opportunities to improve the flexibility, robustness, and accuracy of existing visual servoing methods. Supervised learning-based VS, such as KOVIS [22] and Siame-se(3) [23], leverages a network to extract features for observed and target images, such as key points or latent space features, and subsequently utilizes the intermediate results to train a policy network to predict motion of robotic arm. Siame-se(3) [23] utilized PBVS in simulation for collecting dataset. However, these approaches to opportunistic deep learning mostly focus on the perceptual part and finding the optimal path for multi-view visual servoing is still rarely investigated.

### B. Reinforcement Learning-based VS and Continuous Control

RL offers advantages over deep learning for visual servoing, including adaptability to uncertainty, support for long-term decision-making. However, predicting continuous action spaces using RL is expected to become more challenging due to the exponential growth in the size of continuous action spaces compared to discrete spaces. Lillicrap et al. [24] proposed off-policy algorithm named Deep Deterministic Policy Gradients (DDPG) with models of actor and critic. To address the overestimation problem of DDPG, Fujimoto et al. [25] proposed an algorithm called Twin Delayed Deep Deterministic Policy Gradient (TD3).

To learning optimal VS control policy with less dependency on prior domain knowledge and tackle unexpected disturbances, RL is utilized to train visual servoing control law [4]. Sampedro et al. [4] used DDPG algorithm to build an IBVS controller for multirotor aerial robots. Shi et al. [5] combined Q-learning with fuzzy state coding to adjust the image Jacobian matrix in IBVS efficiently and adaptively, aiming to stabilize and improve the VS performance for wheeled mobile robots. Singh et al. [6] showed that prioritized experience replay buffer could improve the convergence time of RL-based IBVS. Lampe et al. [7] proposed reinforcement learning-based visual servoing for reaching and grasping.

## III. PROBLEM STATEMENT AND METHOD

### A. Problem statement

We formulate the reinforcement learning-based visual servoing as a Markov decision process. The robot agent executes a continuous action $\boldsymbol{a}_t$ estimated by robot control policy, denoted as $\pi_{\mathrm{vs}}$, according to any given state $\boldsymbol{s}_t$ at time $t$. After action execution, reward is calculated based on updated new state $\boldsymbol{s}_{t+1}$. The main objective of RL-based visual servoing is to train an optimal robot control policy $\pi_{\mathrm{vs}}^*$ that maximizes the expected rewards $R$ in the future.

In our work, we investigate a robot visual servoing off-policy to estimate robot joint velocities $\dot{\boldsymbol{q}}_{t+1}$ based on state $\boldsymbol{s}_t$. The policy aims to iteratively achieve the desired pose by maximizing the object function $J$ with respect to state $S$ and keeping robot free of singularities and self-collisions, as demonstrated in Eq. 1.

$$
\begin{aligned}
\boldsymbol{a} &= \pi_{\mathrm{vs}}\left(\boldsymbol{s}; \boldsymbol{w}^p\right) \\
J\left(\boldsymbol{w}^p\right) &= \mathbb{E}_S\left[\widehat{q}\left(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{w}^v\right)\right]
\end{aligned}
\tag{1}
$$

where $\boldsymbol{w}^p$ and $\boldsymbol{w}^v$ denote the weights of policy network (actor) $\pi_{\mathrm{vs}}$ and value network (critic) $\widehat{q}$.

The main architecture of RL-based visual servoing is shown in Fig. 2. Firstly, the latent space representations are learned from visual state, as introduced in Sec. III-B. Then, proposed closed-loop multi-perspective visual servoring network estimated robot actions based on latent space representations and robot states, as detailed in Sec. III-C.

### B. Learning Latent Space Representations from Autoencoder

To acquire generalized features and minimal information of images for visual servoing, an autoencoder network is utilized to extract latent space representations from depth image. The autoencoder is trained with reconstruction loss function based on self-supervison:

$$
\mathcal{L}(\boldsymbol{w}_1, \boldsymbol{w}_2; \boldsymbol{y}) = \mathbb{E}\left[\left(\boldsymbol{y} - \underbrace{\mathrm{P}_{\boldsymbol{w}_1}\left(Q_{\boldsymbol{w}_2}(\boldsymbol{y})\right)}_{\boldsymbol{y}'}\right)^2\right]
\tag{2}
$$

where $\boldsymbol{y}$ represents the network input, network output $\boldsymbol{y}'$ is calculated based on encoder $\mathrm{P}_{\boldsymbol{w}_1}$ and decoder $Q_{\boldsymbol{w}_2}$ with weights $\boldsymbol{w}_1$ and $\boldsymbol{w}_2$.

### C. Closed-Loop Multi-Perspective Visual Servoing Robotic Policy Network

*1) Action Space and Observation Space:* The action space of proposed reinforcement learning-based visual servoing is the velocity of camera frame with minimal and maximal limits, as described as:

$$
\begin{aligned}
\mathcal{A} &= \{{}^c\dot{\boldsymbol{x}}_c \mid {}^c\dot{\boldsymbol{x}}_c \in [{}^c\dot{\boldsymbol{x}}_{c,\min}, {}^c\dot{\boldsymbol{x}}_{c,\max}]\} \\
{}^c\dot{\boldsymbol{x}}_c &= (v_{c,x}, v_{c,y}, v_{c,z}, \omega_{c,x}, \omega_{c,y}, \omega_{c,z})
\end{aligned}
\tag{3}
$$

where $\mathcal{A}$ denotes the action space of RL algorithm. $\dot{\boldsymbol{x}}_c$ represents camera frame velocity, $v_{c,x}, v_{c,y}, v_{c,z}$ and $\omega_{c,x}, \omega_{c,y}, \omega_{c,z}$ denote linear and angular speeds in $x$, $y$, $z$ direction.
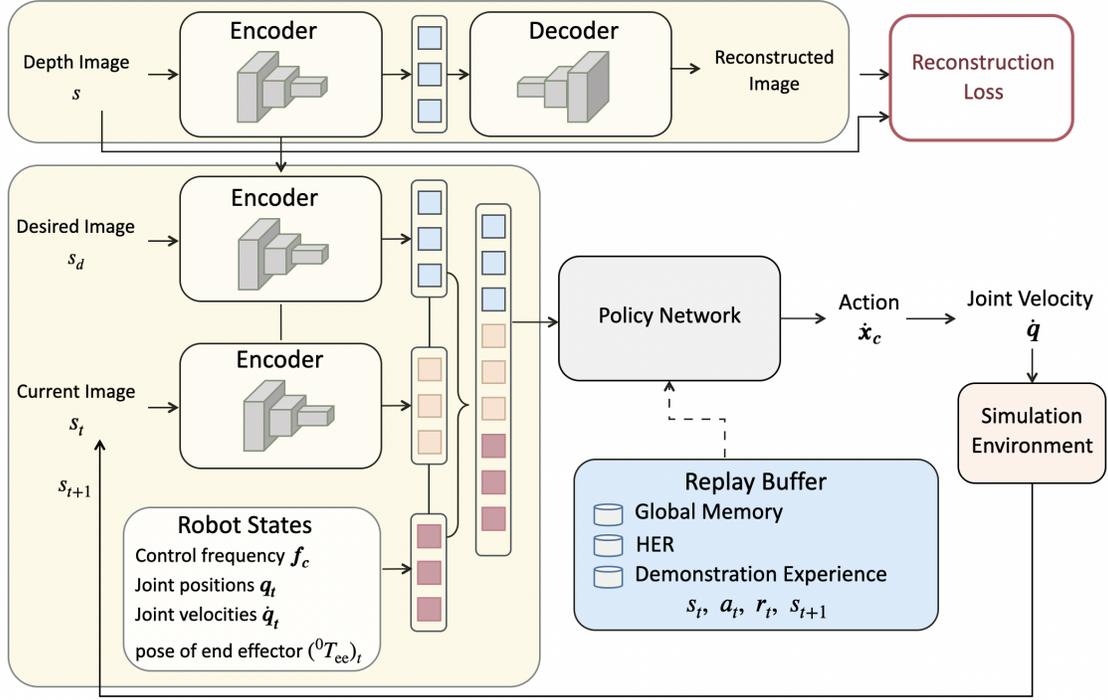
Fig. 2. Architectures of Autoencoder and Closed-Loop Multi-Perspective Visual Servoing Network. Top Part: The latent space representations are extracted with encoder of autoencoder. Bottom Part: The robot policy estimates joint velocities based on latent space visual representations of current and desired image frames and robot states, including joint position, joint velocities, end-effector pose, and control frequency. The action of joint velocities is executed in simulation environment and the policy network is trained using reinforcement learning.

The observation space $\mathcal{S}$ consists of the visual states expressed in latent space $\boldsymbol{S}_t, \boldsymbol{S}_{\text{des}}$ and the robot states $\boldsymbol{S}_{\text{robot}}$:

$$\mathcal{S} = \left\{ \boldsymbol{S}_t, \boldsymbol{S}_{\text{des}}, f_c, \boldsymbol{q}, \dot{\boldsymbol{q}}, {}^{0}\boldsymbol{x}_{\text{ee}} \right\} \quad (4)$$

where $f_c$ denotes control frequency, $\boldsymbol{q}$ and $\dot{\boldsymbol{q}}$ represent joint positions and joint velocities. ${}^{0}\boldsymbol{x}_{\text{ee}}$ denotes the Cartesian pose of end-effector coordinate frame.

*2) Architecture:* Firstly, the real-time captured image $\boldsymbol{I}_t$ at time $t$ and desired image $\boldsymbol{I}_{\text{des}}$ are fed into encoders with shared weights to obtain the latent space features $\boldsymbol{S}_t$ and $\boldsymbol{S}_{\text{des}}$. Secondly, the robot policy network takes the $\boldsymbol{S}_{\text{des}}, \boldsymbol{S}_t$ together with $\boldsymbol{S}_{\text{robot}}$ as inputs. The robot states encompass control frequency $f_c$, joint positions $\boldsymbol{q}$ and velocities $\dot{\boldsymbol{q}}$, as well as the pose ${}^{0}\boldsymbol{T}_{\text{ee}}$ of the end-effector derived from the robot's forward kinematics. The network then produces Cartesian space velocity $\dot{\boldsymbol{x}}_c$ of camera coordinate frame and based on this velocity, we calculate joint velocities $\dot{\boldsymbol{q}}_t$ as follows:

$$\begin{aligned} v_{\text{ee}} &= {}^{\text{ee}}\boldsymbol{V}_c \dot{\boldsymbol{x}}_c \\ \dot{\boldsymbol{q}} &= ({}^{\text{e}}\boldsymbol{J}_{\text{e}})^{-1} v_{\text{ee}} \end{aligned} \quad (5)$$

where ${}^{\text{ee}}\boldsymbol{V}_c$ and ${}^{\text{e}}\boldsymbol{J}_{\text{e}}$ denote twist transformation matrix and Jacobian matrix, $v_{\text{ee}}$ introduces velocity of end effector.

After robotic movement based on joint velocities $\dot{\boldsymbol{q}}_t$ and control frequency $f_c$, the camera image $\boldsymbol{I}_{t+1}$ and robot states $f_c, \boldsymbol{q}_{t+1}, \dot{\boldsymbol{q}}_{t+1}, ({}^{0}\boldsymbol{T}_{\text{ee}})_{t+1}$ are updated and serve as the feedbacks in the closed control visual servoing loop.

*3) Twin-Delayed Deep Deterministic-based Policy Gradient Agents:* We design our robot policy based on the architecture of the actor-critic network TD3. It incorporates a total of six networks, comprising two critic neural networks, one actor neural network, and three target neural networks, as detailed in Fig. 3.

TD3 is proposed based on clipped double Q-learning structure, delayed updates of policy networks, and target networks to solve the overestimation problem. The policy network (actor) estimates the velocities of camera frame and the robot action $\hat{\boldsymbol{a}}_{t+1}$ is executed in our simulation environment.

$$\hat{\boldsymbol{a}}_{t+1} = \pi\left(\boldsymbol{s}_{t+1}; \boldsymbol{w}^p\right) + \epsilon \quad (6)$$

where $\boldsymbol{w}^p$ represents weights of the policy network and $\epsilon$ denotes the noise from a clipped normal distribution.

Based on the output $\hat{q}\left(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{w}^v\right)$ of value network (critic) $Q_\pi\left(\boldsymbol{s}, \boldsymbol{a}; \boldsymbol{w}^v\right)$, temporal difference target (TD-target) $\hat{y}_t$ is calculated as follows:

$$\hat{y}_t = r_t + \eta \hat{q}\left(\boldsymbol{s}_{t+1}, \hat{\boldsymbol{a}}_{t+1}; \boldsymbol{w}^v\right) \quad (7)$$

where $r_t$ denotes reward and $\boldsymbol{w}^v$ represents weights of value network. The value networks are updated based on error of the TD-target $\hat{q}\left(\boldsymbol{s}_t, \boldsymbol{a}_t; \boldsymbol{w}^v\right) - \hat{y}_t$. During updating parameters, the minimal TD-target is selected from two target critics. Taking the minimal TD-target $y$ can not only reduce the overestimation but also mitigate the bias propagation. The update of target policy networks is less frequent than the
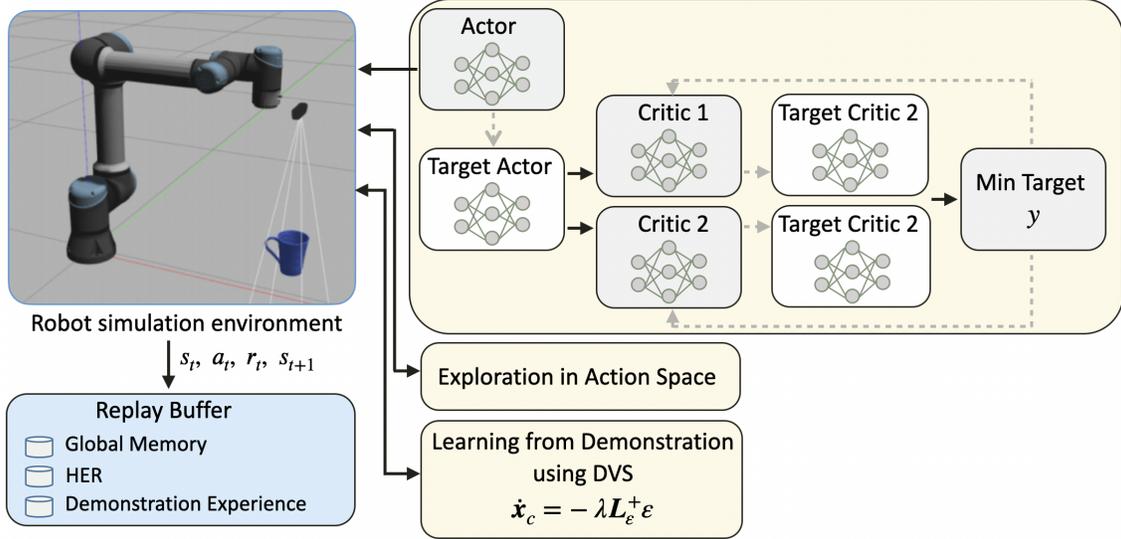
Fig. 3. Architecture of proposed visual servoing policy network. Exploration from action space is utilized to capture more trajectory data. Traditional DVS method is used to generate imperfect success demonstration experience and HER is applied in generating experiences from achieved goals.

update of critic networks to avoid using highly variant value estimates to update the policy. The delayed update improves the performance of policy networks and stabilizes the training process. In each episode, the transition $(s_t, a_t, r_t, s_{t+1})$ will be stored in the replay buffer as a global memory part and sampled as a train set.

*4) Reward function:* We propose reward function $r_t$ considering task- and robot-specific constraints for multi-perspective visual servoing, as formulated in Eq. 8. The reward function considers the translation error $e_{\text{trans}}$, rotational error $e_{\text{rot}}$, and image difference error $e_{\text{img}}$, as well as error related to training step $e_{\text{step}}$ and terminal reward $r_{\text{terminal}}$ with weights $\phi_1, \phi_2, \phi_3, \phi_4$.

$$
\begin{aligned}
r_t =& \phi_1 \left( e_{\text{trans,t}} - e_{\text{trans,t+1}} \right) + \phi_2 \left( e_{\text{rot,t}} - e_{\text{rot,t+1}} \right) + \\
& \phi_3 \left( e_{\text{img,t}} - e_{\text{img,t+1}} \right) - \phi_4 e_{\text{step}} + r_{\text{terminal}},
\end{aligned} \tag{8}
$$

If the robot faces singularity, collision, joint limitation problems, the training will be terminated and transition is stored. The terminal reward is formulated as follows.

$$
r_{\text{terminal}} = \begin{cases} 100 - \|a_t\|, & e_{\text{trans}} < \varphi_{\text{trans}} \wedge \\ & e_{\text{rot}} < \varphi_{\text{rot}} \wedge \neg \text{failure}; \\ -100, & \text{if failure}, \end{cases}
$$

$$
\text{failure} = true, \text{ if} \begin{cases} e_{\text{trans}} > \varphi_{\text{trans}} \vee \\ e_{\text{rot}} > \varphi_{\text{rot}} \vee \\ \det \left( {}^0 J_{ee} \right) \leq \varphi_{\text{Jacobian}} \vee \\ \text{jointlimitsreached} \vee \\ \text{collision detected} \vee \\ \text{object out of FOV} \vee \\ \text{max. step reached} \end{cases} \tag{9}
$$

where $e_{\text{trans,t}}$ and $e_{\text{rot,t}}$ denote translation error and rotational error between camera pose ${}^0T_c$ and ${}^0T_{c^*}$ at time step $t$. $e_{\text{img,t}}$ represents image difference between $I_t$ and $I_{\text{des}}$. $e_{\text{step}}$ introduces constant step error with value 1. The binary

terminal reward of $\pm 100$ is set such that the average episode returns of $+100$ and $-100$ could represent the success and failure respectively. The $\varphi_{\text{trans}}, \varphi_{\text{rot}}, \varphi_{\text{Jacobian}}$ mean thresholds of translation error, rotation error, and Jacobian matrix.

*5) Training Strategy.:* To investigate the efficiency of training RL-based policy for robotic visual servoing with sparse action space, four variants of our proposed method were trained according to following training strategies:

**Hindsight Experience Replay.** To improve the sample efficiency of RL algorithm with sparse reward and sparse data, HER is proposed by learning from failed experiences. Using HER method, achieved states are sampled from failed trajectory as reachable goals and the transitions are stored into replay buffer. Specifically, the policy could learn to achieve an arbitrary given goal from HER experiences.

**Learning from Demonstration.** Due to sparse distribution of success action of visual servoing, it's challenging to learn proposed policy using pure TD3. We propose the method to move the robot agent nearby the desired state and utilize traditional DVS method to collect imperfect demonstration experiences, according to:

$$
\dot{x}_c = -\lambda L_\varepsilon^+ \varepsilon. \tag{10}
$$

where $\dot{x}_c$ denotes camera velocity, $L_\varepsilon^+$ represents Moore-Penrose pseudo-inverse of interaction matrix. $\varepsilon$ denotes visual error and the weight $\lambda$ denotes a positive scalar value.

**Additional Exploration.** We add additional exploration phase before training TD3-based network, where the random actions are sampled before train and replay buffer collects the corresponding transitions with exploration experiences.

Based on above training strategies, four variants are introduced:

**Pure TD3.** Select pure TD3 as baseline method.

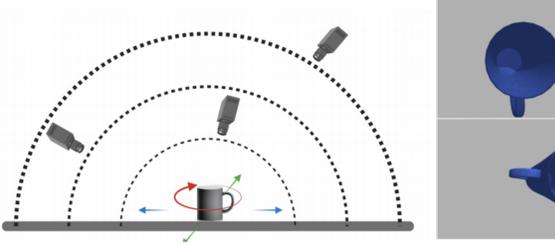| | Setting 1 | Setting 2 | Setting 3 |
|---|---|---|---|
| Range of object movement | ±5 cm | ±10 cm | ±10 cm |
| Random initial pose | No | No | Yes |

Fig. 4. Scene generation for training the autoencoder and rendered images.

**TD3 + Exploration.** The robot policy network is based on the pure TD3 network and additional exploration phase is added before parameter updating of policy and value networks.

**TD3 + HER + exploration.** At the beginning of each episode in both the exploration and the training phase, we generated additional experiences in replay buffer based on HER method.

**TD3 + HER + exploration + Learning from Demonstration.** The imperfect demonstration experiences are generated with a certain probability and stored into replay buffer.

## IV. EXPERIMENTS

We trained proposed closed-loop visual servoing network and executed evaluation experiments in simulation environment. The experiments aimed to 1) investigate reinforcement learning method in multi-perspective visual servoing with sparse data distribution and rewards, 2) compare the traditional visual servoing method and proposed method.

Sec. IV-A describes the experimental setup and pipeline of data collection. Sec. IV-B details the comparison experiments of our proposed methods and baseline methods.

### A. Experimental Setup and Data Collection

The experimental setup was shown in Fig. 3, where UR5e robot agent observed the blue cup on the table with Intel RealSense D435 camera. Before starting visual servoing, the observed object was placed on the table in a range and randomly rotated along z-axis. As well, we designed three experimental settings based on the range of the object and whether the robot's initial pose was randomized or not before visual servoing, as summarized in Tab. I.

To learn visual servoing from latent space representation, we first trained an autoencoder. To achieve this, we generated a dataset where depth and color images are rendered from a
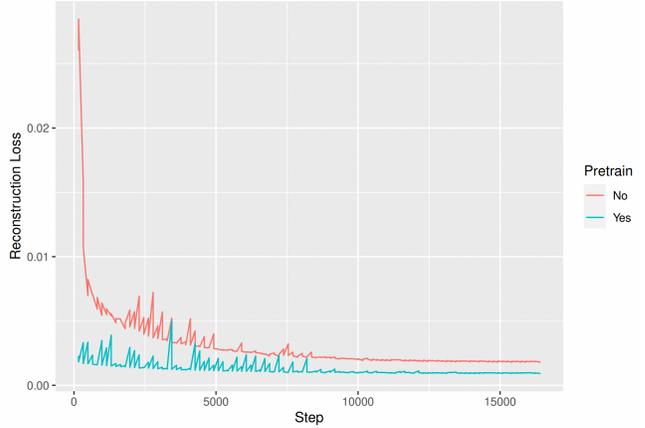


Fig. 5. Training curves of autoencoder with and without pre-trained model.

wide range of different relative poses between the object and camera. As shown in Fig. 4, we generated the random poses of the camera on an upper spherical surface with a random radius. Subsequently, we moved the object pose randomly along the x- and y-axis and rotate it around the z-axis.

During data generation, 100 random poses of the camera were generated and 100 random object poses were synthesized for each camera pose. The radius was selected from $50$mm to $850$mm. Finally, one dataset with 10,000 samples was generated to train the autoencoder to learn the latent space representation of the object from different perspectives. We trained network with pre-training model from CelebA dataset [26] and the training curves with and without pre-trained model were shown in Fig. 5 to demonstrate successful convergence of autoencoder.

After training the autoencoder, the synthetic robot agent was controlled based on the proposed visual servoing method. During the episodes of the training procedure, the simulation environment, value and policy networks, and HER buffer were firstly initialized. The goal state was randomly generated where the observed cup was in the observation area of the synthetic camera. In the training process with maximal step, policy network sampled action $a_t$ and state $s_{t+1}$ were collected after execution of robot agent with velocity control.

### B. Simulation Experiments of Visual Servoing

We executed comparison experiments of the four variants of proposed method, comparison experiments of our proposed method and traditional baseline method DVS.

In comparison experiments of four variants of proposed methods, the first experiment setting was applied during training. The training curves were shown in Fig. 6. The variant TD3+HER+Exploration+Demonstration performed best performance. The pure TD3-based network could not converge for multi-perspective visual servoing from latent space representations. The exploration method slightly improved the performance of the pure TD3-based variant. HER-based variants can converge successfully. This suggests that HER is beneficial in mitigating the problem of sparse data and
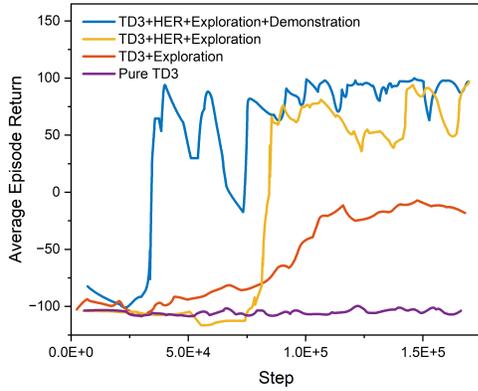
Fig. 6. Training curves of four variants in first experimental configuration.



Fig. 7. Success rates of DVS and our method in different experimental settings. The vacant bar within the bar chart symbolizes a success rate of 0%.

reward function. The variant using learning from demonstration could converge faster.

In comparison experiments of our method and traditional baseline method, we conducted 100 comparison experiments under each experimental setting. We chose DVS as the traditional baseline method due to the lack of distinctive visual features on the observed object.

Evaluation metrics were considered, including success rate and error distribution in translation and rotation. The visual servoing was counted as success when the error was smaller than a threshold. During training, we set the threshold as 2mm. Secondly, the translation and rotation errors were summarized based on a series of experiments in simulation.

Our method significantly outperformed DVS in terms of success rate. The mean success rate of the proposed method achieved above $97.0\%$ in different experimental settings, as shown in Fig. 7. Meanwhile, the traditional visual servoing method DVS performed a poor success rate in complicated experimental settings. The mean success rate of experimental settings achieved $42.3\%$.

We also quantitatively evaluated with the error distributions in translation and rotation of multi-perspective visual servoing in three experimental settings, as depicted in Fig. 8. The cases of failed convergence to the target state were shown as empty bars. The error threshold of 2 mm was utilized during training our model. In simple scenarios of setting 1 and setting 2, we observed a slight increase in errors in proposed RL-based algorithm compared to DVS method. Our approach involves a trade-off between exploration and exploitation. Specifically, in simple scenarios, where problems are relatively straightforward, traditional algorithms may tend to exploit known information, while our reinforcement learning-based algorithm may lean towards exploring new possibilities. This tendency might lead to a minor increase in errors in simple scenarios. However, in complex scenarios, this exploratory strategy could be more advantageous in finding better solutions. Our algorithm demonstrates greater adaptability. It can learn and adjust in
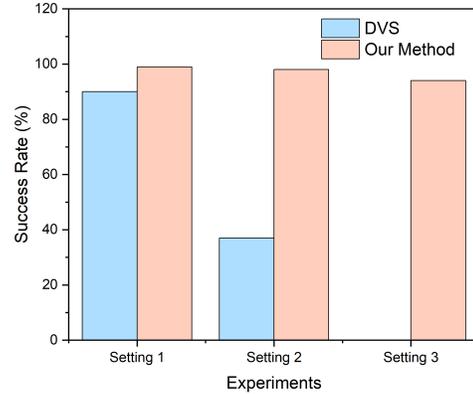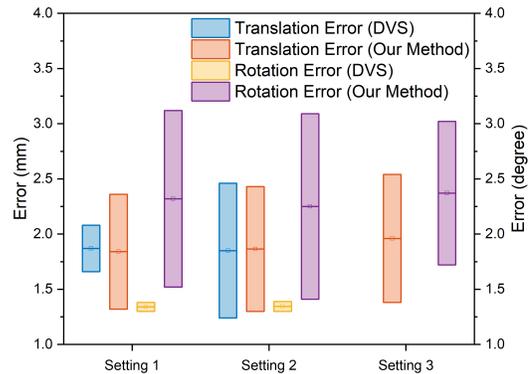


Fig. 8. Translation and rotational error distributions of DVS and our method (error threshold = 2 mm) in different configuration settings.
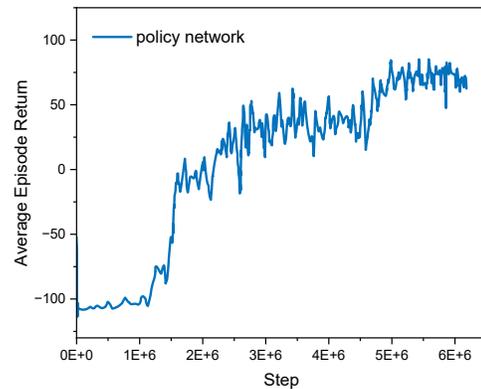


Fig. 9. Training curve with average return per episode of proposed method in third experimental configuration.

different scenarios. In complex situations, where there is significant environmental variation and uncertainty, traditional algorithms might struggle to adapt. In contrast, our algorithm excels in handling complexity and performs better in such scenarios. In summary, although our method may experience a slight increase in errors in simple scenarios, it exhibits outstanding performance in complex scenarios. This suggests that our approach is more robust and adaptable, enabling it to function effectively in the intricate environments of the real world.

Finally, we trained proposed method using HER, learning from demonstration, additional exploration in third experimental setting with more complicated scenes. The training took 141 hours, as demonstrated in Fig. 9.

## V. CONCLUSIONS AND FUTURE WORK

We proposed a novel closed-loop multi-perspective reinforcement learning-based visual servoing network. HER, learning from demonstration and additional exploration methods were utilized to alleviate the pain of convergence in sparse reward function and sparse success behaviors in action space. The robot agent with velocity control was developed in simulation to perform visual servoing with different complicated scenarios. The robot actions were estimated based on latent space representations learned from visual states.

The comparison experiments proved that our variant TD3+HER+Exploration+Demonstration demonstrated ability of our method in multi-perspective visual servoing using reinforcement learning. From quantitative experiments of our method and traditional method DVS, our method outperformed the DVS with a mean success rate of $97.0\%$ in different experimental settings. Meanwhile, our method could converge in complicated scenes with desired error distributions in translation and rotation. The mean translation errors of our method achieved performances of the success cases of accurate traditional method DVS. The future work will extend the reinforcement learning network in robot arm and five-finger hand manipulation.

## REFERENCES

[1] A. Saxena, H. Pandya, G. Kumar, A. Gaud, and K. M. Krishna, "Exploring convolutional networks for end-to-end visual servoing," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 3817–3823.

[2] Q. Bateux, E. Marchand, J. Leitner, F. Chaumette, and P. Corke, "Training deep neural networks for visual servoing," in *2018 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2018, pp. 3307–3314.

[3] C. Yu, Z. Cai, H. Pham, and Q.-C. Pham, "Siamese convolutional neural network for sub-millimeter-accurate camera pose estimation and visual servoing," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2019, pp. 935–941.

[4] C. Sampedro, A. Rodriguez-Ramos, I. Gil, L. Mejias, and P. Campoy, "Image-based visual servoing controller for multirotor aerial robots using deep reinforcement learning," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 979–986.

[5] H. Shi, H. Wu, C. Xu, J. Zhu, M. Hwang, and K.-S. Hwang, "Adaptive image-based visual servoing using reinforcement learning with fuzzy state coding," *IEEE Transactions on Fuzzy Systems*, vol. 28, no. 12, pp. 3244–3255, 2020.

[6] P. Singh, V. Singh, S. Dutta, and S. Kumar, "Model & feature agnostic eye-in-hand visual servoing using deep reinforcement learning with prioritized experience replay," in *2019 28th IEEE International Conference on Robot and Human Interactive Communication (RO-MAN)*. IEEE, 2019, pp. 1–8.

[7] T. Lampe and M. Riedmiller, "Acquiring visual servoing reaching and grasping skills using neural reinforcement learning," in *The 2013 international joint conference on neural networks (IJCNN)*. IEEE, 2013, pp. 1–8.

[8] H. Shi, X. Li, K.-S. Hwang, W. Pan, and G. Xu, "Decoupled visual servoing with fuzzy q-learning," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 1, pp. 241–252, 2016.

[9] Z. Jin, J. Wu, A. Liu, W.-A. Zhang, and L. Yu, "Policy-based deep reinforcement learning for visual servoing control of mobile robots with visibility constraints," *IEEE Transactions on Industrial Electronics*, vol. 69, no. 2, pp. 1898–1908, 2021.

[10] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, O. Pieter Abbeel, and W. Zaremba, "Hindsight experience replay," *Advances in neural information processing systems*, vol. 30, 2017.

[11] L. Weiss, A. Sanderson, and C. Neuman, "Dynamic sensor-based control of robots with visual feedback," *IEEE Journal on Robotics and Automation*, vol. 3, no. 5, pp. 404–417, 1987.

[12] J. T. Feddema and O. R. Mitchell, "Vision-guided servoing with feature-based trajectory generation (for robots)," *IEEE Transactions on Robotics and Automation*, vol. 5, no. 5, pp. 691–700, 1989.

[13] W. J. Wilson, C. W. Hulls, and G. S. Bell, "Relative end-effector control using cartesian position based visual servoing," *IEEE Transactions on Robotics and Automation*, vol. 12, no. 5, pp. 684–696, 1996.

[14] B. Thuilot, P. Martinet, L. Cordesses, and J. Gallice, "Position based visual servoing: keeping the object in the field of vision," in *Proceedings 2002 IEEE International Conference on Robotics and Automation (Cat. No. 02CH37292)*, vol. 2. IEEE, 2002, pp. 1624–1629.

[15] F. Chaumette, "Potential problems of stability and convergence in image-based and position-based visual servoing," in *The confluence of vision and control*. Springer, 1998, pp. 66–78.

[16] P. I. Corke and S. A. Hutchinson, "A new partitioned approach to image-based visual servo control," *IEEE Transactions on Robotics and Automation*, vol. 17, no. 4, pp. 507–515, 2001.

[17] K. Deguchi, "A direct interpretation of dynamic images with camera and object motions for vision guided robot control," *International Journal of Computer Vision*, vol. 37, no. 1, pp. 7–20, 2000.

[18] V. Kallem, M. Dewan, J. P. Swensen, G. Hager, and N. J. Cowan, "Kernel-based visual servoing," *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 1975–1980, 2007.

[19] C. Collewet, E. Marchand, and F. Chaumette, "Visual servoing set free from image processing," in *2008 IEEE International Conference on Robotics and Automation*. IEEE, 2008, pp. 81–86.

[20] A. Dame and E. Marchand, "Entropy-based visual servoing," in *2009 IEEE International Conference on Robotics and Automation*. IEEE, 2009, pp. 707–713.

[21] C. Collewet and E. Marchand, "Photometric visual servoing," *IEEE Transactions on Robotics*, vol. 27, no. 4, pp. 828–834, 2011.

[22] E. Y. Puang, K. P. Tee, and W. Jing, "Kovis: Keypoint-based visual servoing with zero-shot sim-to-real transfer for robotics manipulation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020, pp. 7527–7533.

[23] S. Felton, E. Fromont, and E. Marchand, "Siame-se (3): regression in se (3) for end-to-end visual servoing," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 14 454–14 460.

[24] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," *arXiv preprint arXiv:1509.02971*, 2015.

[25] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. PMLR, 10–15 Jul 2018, pp. 1587–1596. [Online]. Available: https://proceedings.mlr.press/v80/fujimoto18a.html

[26] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 3730–3738.