# MPC-based Reinforcement Learning for a Simplified Freight Mission of Autonomous Surface Vehicles

Wenqi Cai, Arash B. Kordabad, Hossein N. Esfahani, Anastasios M. Lekkas, and Sébastien Gros

*Abstract*—In this work, we propose a Model Predictive Control (MPC)-based Reinforcement Learning (RL) method for Autonomous Surface Vehicles (ASVs). The objective is to find an optimal policy that minimizes the closed-loop performance of a simplified freight mission, including collision-free path following, autonomous docking, and a skillful transition between them. We use a parametrized MPC-scheme to approximate the optimal policy, which considers path-following/docking costs and states (position, velocity)/inputs (thruster force, angle) constraints. The Least Squares Temporal Difference (LSTD)-based Deterministic Policy Gradient (DPG) method is then applied to update the policy parameters. Our simulation results demonstrate that the proposed MPC-LSTD-based DPG method could improve the closed-loop performance during learning for the freight mission problem of ASV.

## I. INTRODUCTION

Autonomous Surface Vehicles (ASVs) are widely applied for many fields, such as freight transportation, military, search and rescue [1], and therefore attract broad attention for scientific and industrial researches. Various methods have been proposed to solve the problem of operating and automating the ASV, including path following, collision avoidance, and autonomous docking [2], [3]. However, designing a control strategy that could realize both collision-free path following and docking in a freight mission with time-varying disturbances is still a topic worth exploring. With the development of Machine Learning (ML), Reinforcement Learning (RL)-based control strategies are getting noticed by people, as they can exploit real data to reduce the impact of model uncertainties and disturbances.

Deterministic Policy Gradient (DPG), as the direct RL method, estimates the optimal policy by a parameterized function approximator, and optimizes the policy parameters directly via gradient descent steps of the performance [4]. Deep Neural Networks (DNNs) are very commonly used function approximators in RL [5]. However, DNN-based RL lacks the abilities concerning the closed-loop stability analysis, state/input constraints satisfaction, and meaningful weights initialization [6]. To address these problems, the perspective of using Model Predictive Control (MPC)-based RL has been proposed and justified in [7], i.e. it suggests using MPC as the function approximation for the optimal policy in RL. Unlike DNNs, MPC-based policies satisfy the state/input constraints and safety requirements by construction, and its well-structured property enables the stability analysis of the system.

However, for computational reasons, simple models are usually preferred in the MPC-scheme. Hence, the MPC model

The authors are with Department of Engineering Cybernetics, Norwegian University of Science and Technology (NTNU), Trondheim, Norway. E-mail:{wenqi.cai,arash.b.kordabad,hossein.n.esfahani,anastasios.lekkas,sebastien.gros}@ntnu.no

often does not have the required structure to correctly capture the real system dynamics and stochasticity. As a result, MPC can deliver a reasonable approximation of the optimal policy, but it is usually suboptimal [8]. Besides, choosing the model parameters that best fit the MPC model to the real system does not necessarily yield a policy that achieves the best closed-loop performance [6]. Therefore, choosing appropriate MPC parameters to achieve the best closed-loop performance is extremely challenging. Nevertheless, according to *Theorem 1* and *Corollary 2* in [7], it can conclude that by adjusting not only the MPC model parameters but also the parameters in the MPC cost and constraints, the MPC scheme can, theoretically, generate the optimal closed-loop policy even if the MPC model is inaccurate. It is also shown that RL is a suitable candidate to perform that adjustment. Recent researches focused on the MPC-based RL have further developed this approach [9], [10], [11], [12], [13].

The contribution of this work is to provide a promising approach for a complete ASV freight mission problem. The problem is challenging since it needs to solve the obstacle avoidance, path following, and autonomous docking simultaneously, in a stochastic environment. We elaborate the proposed MPC-based RL method in the ASV problem framework, as well as formulate an algorithm for the MPC-LSTD-based DPG method.

## II. ASV MODEL

The 3-Degree of Freedom (3-DOF) position of the vessel can be represented by a pose vector $\boldsymbol{\eta} = [x, y, \psi]^T \in \mathbb{R}^3$ in the North-East-Down (NED) frame, where $x$ is the North position, $y$ is the East position, and $\psi$ is the heading angle (see Fig. 1). The velocity vector $\boldsymbol{\nu} = [u, v, r]^T \in \mathbb{R}^3$, including the surge velocity $u$, sway velocity $v$, and yaw rate $r$, is decomposed in the body-fixed frame. The nonlinear dynamics can be written
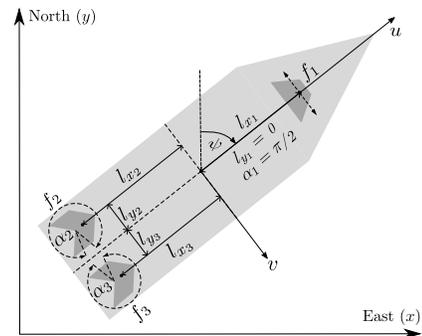


Fig. 1. The 3-DOF ASV model in the NED frame.

as follows [14]

$$\dot{\boldsymbol{\eta}} = \boldsymbol{J}(\psi)\boldsymbol{\nu} \tag{1a}$$

$$\boldsymbol{M}\dot{\boldsymbol{\nu}} + \boldsymbol{D}\boldsymbol{\nu} = \boldsymbol{\tau} + \boldsymbol{\tau}_a, \tag{1b}$$

where $\boldsymbol{J}(\psi) \in \mathbb{R}^{3\times3}$ is the rotation matrix, $\boldsymbol{M} \in \mathbb{R}^{3\times3}$ is the mass matrix, and $\boldsymbol{D} \in \mathbb{R}^{3\times3}$ is the damping matrix (see [15] for their specific physical meanings and values). Vector $\boldsymbol{\tau} \in \mathbb{R}^3$ presents the control forces and moment empowered by the thrusters. Vector $\boldsymbol{\tau}_a \in \mathbb{R}^3$ is the additional forces rendered from disturbances, e.g., wind, ocean wave and etc. The thrust configuration is illustrated in Fig. 1. The vector $\boldsymbol{\tau}$ could be specifically written as $\boldsymbol{\tau} = \boldsymbol{T}(\boldsymbol{\alpha})\boldsymbol{f}$, where $\boldsymbol{f} = [f_1, f_2, f_3]^\top \in \mathbb{R}^3$ is the thruster forces vector as we consider one tunnel thruster $f_1$ and two azimuth thrusters $f_2, f_3$. They are subjected to the bounds

$$f_{p\min} \le f_p \le f_{p\max}, \quad p = 1, 2, 3. \tag{2}$$

Matrix $\boldsymbol{T}(\boldsymbol{\alpha}) \in \mathbb{R}^{3\times3}$ presents the thruster configuration, written as

$$\boldsymbol{T}(\boldsymbol{\alpha}) = \begin{bmatrix} 0 & \cos(\alpha_2) & \cos(\alpha_3) \\ 1 & \sin(\alpha_2) & \sin(\alpha_3) \\ l_{x1} & T_{32} & T_{33} \end{bmatrix}, \tag{3}$$

where elements $T_{32} = l_{x2}\sin(\alpha_2) - l_{y2}\cos(\alpha_2)$, and $T_{33} = l_{x3}\sin(\alpha_3) - l_{y3}\cos(\alpha_3)$. Constants $l_{xi}$ and $l_{yi}$ with $i = 1, 2, 3$ are the distances between each thruster and the cross line of the ship's center. Term $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \alpha_3]^\top \in \mathbb{R}^3$ is the corresponding orientation vector. The angle $\alpha_1$ is fixed ($\pi/2$), while $\alpha_2$ and $\alpha_3$, associated to the two azimuth thrusters, are restricted in the range

$$|\alpha_2 + \pi/2| \le \alpha_{\max}, \quad |\alpha_3 - \pi/2| \le \alpha_{\max}. \tag{4}$$

A maximum angle of $\alpha_{\max}$ with a forbidden sector is considered in this work to avoid thrusters 2 and 3 directly work against each other, as shown in Fig. 1. With a sampling time of $dt$, we discretize the ship system (1) as

$$\boldsymbol{s}_{k+1} = F(\boldsymbol{s}_k, \boldsymbol{a}_k, \boldsymbol{\tau}_a), \tag{5}$$

where $\boldsymbol{s}_k = [\boldsymbol{\eta}_k^\top, \boldsymbol{\nu}_k^\top]^\top$ and $\boldsymbol{a}_k = [\boldsymbol{f}_k^\top, \boldsymbol{\alpha}_k^\top]^\top$ are system state and input vectors, respectively. Subscript $k$ denotes the physical time and $F(\cdot)$ is the discretized real system.

## III. SIMPLIFIED FREIGHT MISSION

In this work, we consider a simplified freight mission problem: the ASV starts from an origin $\boldsymbol{A}$ to the end $\boldsymbol{B}$, which is supposed to follow a designed collision-free course and finally dock at the wharf autonomously. Note that the transition from path following to docking is a notable point of this problem.

### A. Collision-Free Path Following

Given a reference path $P_{\text{ref}}$. At time instance $k$, $\boldsymbol{P}_k^{\text{ref}} = [x_k^{\text{ref}}, y_k^{\text{ref}}]^\top$. Then path following could be thought as minimizing the error $l(\boldsymbol{\eta}_k)$

$$l(\boldsymbol{\eta}_k) = \left\| \boldsymbol{\eta}_k^p - \boldsymbol{P}_k^{\text{ref}} \right\|_2^2 = (x_k - x_k^{\text{ref}})^2 + (y_k - y_k^{\text{ref}})^2, \tag{6}$$

where $\boldsymbol{\eta}_k^p = [x_k, y_k]^\top$ contains the first two elements of $\boldsymbol{\eta}_k$. Besides, we assume obstacles of round shape. To avoid these obstacles, the following term $g_n(\boldsymbol{\eta}_k)$, representing the position of the ship relative to the $n^{\text{th}}$ obstacle, should satisfy

$$(x_k - o_{x,n})^2 + (y_k - o_{y,n})^2 \ge (r_n + r_o)^2, \tag{7}$$

i.e.,

$$\underbrace{1 - \left( (x_k - o_{x,n})^2 + (y_k - o_{y,n})^2 \right) \big/ (r_n + r_o)^2 \le 0,}_{g_n(\boldsymbol{\eta}_k)} \tag{8}$$

where $(o_{x,n}, o_{y,n})$ and $r_n$ are the center and radius of the $n^{\text{th}}$ circular obstacle ($n = 1, \dots, N_o$), respectively. Constant $r_o$ is the radius of the vessel and $N_o$ is the number of obstacles.

### B. Autonomous Docking

Docking refers to stopping the vessel exactly at the endpoint $\boldsymbol{B}$ as well as avoiding collisions between any part of the vessel and the quay [3]. The "accurate stop" requires not only an accurate docking position but also zero-valued velocities and thruster forces at the final time, i.e., we ought to minimize

$$h(\boldsymbol{\eta}_k, \boldsymbol{\nu}_k, \boldsymbol{f}_k) = \|\boldsymbol{\eta}_k - \boldsymbol{\eta}_d\|_2^2 + \|\boldsymbol{\nu}_k\|_2^2 + \|\boldsymbol{f}_k\|_2^2, \tag{9}$$

where $\boldsymbol{\eta}_d = (x_d, y_d, \psi_d)$ is the desired docking position. Successfully docking requires $h(\boldsymbol{\eta}_K, \boldsymbol{\nu}_K, \boldsymbol{f}_K) \approx 0$, where subscript $K$ denotes the terminal time step of the freight mission. As for "collision avoidance", we define a safety operation region $\mathbb{S}$ as the spatial constraints for the vessel. The operation region is chosen as the largest convex region that encompasses the docking point but not intersecting with the land. Thus, as long as the vessel is within the region $\mathbb{S}$, no collision will occur during docking, i.e. the following condition should hold

$$\boldsymbol{\eta}_k^p \in \mathbb{S}, \quad \mathbb{S} = \{\boldsymbol{x} | \boldsymbol{A}\boldsymbol{x} < \boldsymbol{b}\}, \tag{10}$$

where $\boldsymbol{\eta}_k^p = [x_k, y_k]^\top$ describes the position of the vessel. The matrix $\boldsymbol{A}$ and the vector $\boldsymbol{b}$ are determined by the shape of the quay and together define the convex region $\mathbb{S}$.

### C. Objective Function

In the context of RL, we seek a control policy $\boldsymbol{\pi}$ that minimizes the following closed-loop performance $J$

$$J(\boldsymbol{\pi}) = \mathbb{E}_{\boldsymbol{\pi}} \left[ \sum_{k=0}^{K} \gamma^k L(\boldsymbol{s}_k, \boldsymbol{a}_k) \Big| \boldsymbol{a}_k = \boldsymbol{\pi}(\boldsymbol{s}_k) \right], \tag{11}$$

where $\gamma \in (0, 1]$ is the discount factor. Expectation $\mathbb{E}_{\boldsymbol{\pi}}$ is taken over the distribution of the Markov chain in the closed-loop under policy $\boldsymbol{\pi}$. The RL-stage cost $L(\boldsymbol{s}_k, \boldsymbol{a}_k)$, in this problem, is defined as a piecewise function:

$$L = \begin{cases} l(\boldsymbol{\eta}_k) + O(\boldsymbol{\eta}_k) + \xi(\boldsymbol{\alpha}_k) & \|\boldsymbol{\eta}_k - \boldsymbol{\eta}_d\|_2^2 > d \\ h(\boldsymbol{\eta}_k, \boldsymbol{\nu}_k, \boldsymbol{f}_k) + \Gamma(\boldsymbol{\eta}_k) + \xi(\boldsymbol{\alpha}_k) & \|\boldsymbol{\eta}_k - \boldsymbol{\eta}_d\|_2^2 \le d, \end{cases} \tag{12}$$

where $O(\boldsymbol{\eta}_k)$ is the obstacle penalty for path following

$$O(\boldsymbol{\eta}_k) = \sum_{n=1}^{N_o} c_n \cdot \max(0, g_n(\boldsymbol{\eta}_k) + d_s), \tag{13}$$

where $c_n > 0$ is the penalty weight, constant $d_s > 0$ is the desired safe distance between vessel and obstacles. Therefore, once the ship breaks the safe distance, i.e. $g_n(\boldsymbol{\eta}_k) + d_s > 0$, a positive penalty will be introduced to the objective function. Function $\Gamma(\boldsymbol{\eta}_k)$ is the collision penalty for docking

$$\Gamma(\boldsymbol{\eta}_k) = \kappa \cdot (1 - \mathbf{1}_{\mathbb{S}}(\boldsymbol{\eta}_k^p)), \tag{14}$$

where $\kappa > 0$ is the penalty weight and $\mathbf{1}_{\mathbb{S}}(\cdot)$ is the indicator function. When the ship is out of the safe region, i.e. $\boldsymbol{\eta}_k^p \notin \mathbb{S}$, a positive penalty will be imposed in the objective function. Function $\xi(\boldsymbol{\alpha}_k)$ is the singular configuration penalty, aiming to avoid the thruster configuration matrix $\boldsymbol{T}(\boldsymbol{\alpha}_k)$ in (3) being singular [16]

$$\xi(\boldsymbol{\alpha}_k) = \frac{\rho}{\varepsilon + \det\left(\boldsymbol{T}(\boldsymbol{\alpha}_k)\boldsymbol{W}^{-1}\boldsymbol{T}^\top(\boldsymbol{\alpha}_k)\right)}, \tag{15}$$

where "det" stands for the determinant of the matrix. Constant $\varepsilon > 0$ is a small number to avoid division by zero, $\rho > 0$ is the weighting of maneuverability, and $\boldsymbol{W}$ is a diagonal weighting matrix. Constant $d > 0$ is designed to substitute the stage cost from path following to docking at $\|\boldsymbol{\eta}_k - \boldsymbol{\eta}_d\|_2^2 = d$, which means that our target transits from path-following to docking when the ship approaches the destination.

## IV. MPC-BASED REINFORCEMENT LEARNING

The core idea of our proposed approach is to use a parameterized MPC-scheme as the policy approximation function, and apply the LSTD-based DPG method to update the parameters so as to improve the closed-loop performance.

### A. MPC-Based Policy Approximation

Consider the following MPC-scheme parameterized with $\boldsymbol{\theta}$

$$\min_{\hat{\boldsymbol{\eta}},\hat{\boldsymbol{\nu}},\hat{\boldsymbol{f}},\hat{\boldsymbol{\alpha}},\boldsymbol{\sigma}} \frac{\theta_d}{\|\hat{\boldsymbol{\eta}}_N - \boldsymbol{\eta}_d\|_2^2 + \delta} \cdot \left(h_{\boldsymbol{\theta}}(\hat{\boldsymbol{\eta}}_N, \hat{\boldsymbol{\nu}}_N) + \Gamma_{\boldsymbol{\theta}}(\hat{\boldsymbol{\eta}}_N)\right) +$$

$$\boldsymbol{\omega}_f^\top \boldsymbol{\sigma}_N + \sum_{i=0}^{N-1} \gamma^i \left(l_{\boldsymbol{\theta}}(\hat{\boldsymbol{\eta}}_i) + \xi(\hat{\boldsymbol{\alpha}}_i) + \boldsymbol{\omega}^\top \boldsymbol{\sigma}_i\right) \tag{16a}$$

$$\text{s.t.} \quad \forall i = 0, \ldots, N-1, \; n = 1, \ldots, N_o$$

$$\left[\hat{\boldsymbol{\eta}}^\top_{i+1}, \hat{\boldsymbol{\nu}}^\top_{i+1}\right]^\top = F_{\boldsymbol{\theta}}(\hat{\boldsymbol{\eta}}_i, \hat{\boldsymbol{\nu}}_i, \hat{\boldsymbol{f}}_i, \hat{\boldsymbol{\alpha}}_i, \boldsymbol{\theta}_a) \tag{16b}$$

$$f_{p_{\min}} \le \hat{f}_{p,i} \le f_{p_{\max}}, \quad p = 1, 2, 3 \tag{16c}$$

$$|\hat{\alpha}_{2,i} + \pi/2| \le \alpha_{\max}, \; |\hat{\alpha}_{3,i} - \pi/2| \le \alpha_{\max}, \tag{16d}$$

$$g_n(\hat{\boldsymbol{\eta}}_i) + \theta_g \le \sigma_{n,i}, \; g_n(\hat{\boldsymbol{\eta}}_N) + \theta_g \le \sigma_{n,N}, \tag{16e}$$

$$\boldsymbol{\sigma}_i \ge 0, \; \boldsymbol{\sigma}_N \ge 0, \tag{16f}$$

$$\hat{\boldsymbol{\eta}}_0 = \boldsymbol{\eta}_k, \; \hat{\boldsymbol{\nu}}_0 = \boldsymbol{\nu}_k, \tag{16g}$$

where $N$ is the prediction horizon. Arguments $\hat{\boldsymbol{\eta}} = \{\hat{\boldsymbol{\eta}}_0, \ldots, \hat{\boldsymbol{\eta}}_N\}$, $\hat{\boldsymbol{\nu}} = \{\hat{\boldsymbol{\nu}}_0, \ldots, \hat{\boldsymbol{\nu}}_N\}$, $\hat{\boldsymbol{f}} = \{\hat{\boldsymbol{f}}_0, \ldots, \hat{\boldsymbol{f}}_{N-1}\}$, $\hat{\boldsymbol{\alpha}} = \{\hat{\boldsymbol{\alpha}}_0, \ldots, \hat{\boldsymbol{\alpha}}_{N-1}\}$, and $\boldsymbol{\sigma} = \{\boldsymbol{\sigma}_0, \ldots, \boldsymbol{\sigma}_N\}$ are the primal decision variables. The term $\frac{\theta_d}{\|\hat{\boldsymbol{\eta}}_N - \boldsymbol{\eta}_d\|_2^2 + \delta} \cdot (h_{\boldsymbol{\theta}}(\cdot) + \Gamma_{\boldsymbol{\theta}}(\cdot))$ introduces a gradually increasing terminal cost as the ship approaches the endpoint, where $\delta > 0$ is a small constant to avoid division by zero. The weighting parameter $\theta_d$, designed to balance the priority of path following and docking, is tuned by RL. Note that $\theta_d$ is chosen to minimize the closed-loop performance that considering both path following and docking, although it may be suboptimal for either single problem.

Parameter $\theta_g$ is the tightening variable used to adjust the strength of the collision avoidance constraints. If the value of $\theta_g$ (positive) is larger, it means that the constraints are tighter and the ship is supposed to be farther away from the obstacles. It is important to use RL to pick an appropriate $\theta_g$, since when $\theta_g$ is too large, although we ensure that the ship safely avoids obstacles, the path following error is increased. Conversely, a smaller $\theta_g$ reduces the following error, but we may gain more penalty when the vessel breaks the safe distance, as described in (13). Note that the obstacle penalties are considered directly as constraints (16e) in the MPC rather than as penalties in the MPC cost, because (8) is a conservative model of the obstacle penalty (13). Variables $\boldsymbol{\sigma}_i$ ($\boldsymbol{\sigma}_i = \{\sigma_{1,i}, \ldots, \sigma_{N_0,i}\}$) and $\boldsymbol{\sigma}_N$ ($\boldsymbol{\sigma}_N = \{\sigma_{1,N}, \ldots, \sigma_{N_0,N}\}$) are slacks for the relaxation of the state constraints, weighted by the positive vectors $\boldsymbol{\omega}$ and $\boldsymbol{\omega}_f$. The relaxation prevents the infeasibility of the MPC in the presence of some hard constraints.

The parameterized stage cost $l_{\boldsymbol{\theta}}(\cdot)$, terminal cost $h_{\boldsymbol{\theta}}(\cdot)$, and docking collision penalty $\Gamma_{\boldsymbol{\theta}}(\cdot)$ in the MPC cost (16a) are designed as follows

$$l_{\boldsymbol{\theta}} = \left\|\hat{\boldsymbol{\eta}}_i^p - \boldsymbol{P}_i^{\text{ref}}\right\|_{\boldsymbol{\Theta}_l}^2 \tag{17a}$$

$$h_{\boldsymbol{\theta}} = \|\hat{\boldsymbol{\eta}}_N - \boldsymbol{\eta}_d\|_{\boldsymbol{\Theta}_\eta}^2 + \|\hat{\boldsymbol{\nu}}_N\|_{\boldsymbol{\Theta}_\nu}^2 \tag{17b}$$

$$\Gamma_{\boldsymbol{\theta}} = \theta_\kappa \cdot (1 - \mathbf{1}_{\mathbb{S}}(\hat{\boldsymbol{\eta}}_N^p)), \tag{17c}$$

where $\boldsymbol{\Theta}_l, \boldsymbol{\Theta}_\eta, \boldsymbol{\Theta}_\nu \in \mathbb{R}^{3\times3}$ are the weighing matrices that are symmetric semi-positive definite. They are expressed as $\boldsymbol{\Theta}_l = (\text{diag}(\boldsymbol{\theta}_l))^2$, $\boldsymbol{\Theta}_\eta = (\text{diag}(\boldsymbol{\theta}_\eta))^2$, $\boldsymbol{\Theta}_\nu = (\text{diag}(\boldsymbol{\theta}_\nu))^2$. Operator "diag" assigns the vector elements onto the diagonal elements of a square matrix. Parameter $\theta_\kappa$ is treated as a degree of freedom for the docking collision penalty. The real model is (5) and we assume the disturbance $\boldsymbol{\tau}_a$ follows a Gaussian distribution. To address the disturbance without using a complex stochastic model in the MPC scheme, one measure is to use a parameter vector $\boldsymbol{\theta}_a \in \mathbb{R}^3$ to parameterize the model as $F_{\boldsymbol{\theta}}(\hat{\boldsymbol{s}}_i, \hat{\boldsymbol{a}}_i, \boldsymbol{\theta}_a)$. As detailed in [7], the full adaptation of the parametrized MPC scheme (model, costs, constraints) can compensate for that unmodelled disturbance. Overall, the adjustable parameters vector $\boldsymbol{\theta}$ is consisted as

$$\boldsymbol{\theta} = \{\boldsymbol{\theta}_l, \boldsymbol{\theta}_\eta, \boldsymbol{\theta}_\nu, \boldsymbol{\theta}_a, \theta_\kappa, \theta_d, \theta_g\}. \tag{18}$$

And $\boldsymbol{\theta}$ will be adjusted by RL according to the principle of "improving the closed-loop performance". Note that: 1. the span of the RL ($K \approx 550$) is much longer than the horizon of the MPC ($N = 60$); 2. the RL cost (12) is a "switching" function, while the MPC cost (16a) contains simultaneously the path following and docking cost to avoid the mixed-integer treatment of the problem; 3. the MPC model does not perfectly match the real system. For the above reasons, having different cost functions in the MPC scheme and RL is rational [7]. Therefore, in order to improve the closed-loop performance of the MPC scheme as assessed by the RL cost, it can be beneficial to parameterize the MPC cost functions, model, and constraints. RL then adjusts these parameters according to the principle of "improving the closed-loop performance". From *Theorem 1* and *Corollary 2* in [7], we know that, theoretically, under some assumptions, if the parametrization is rich enough,

the MPC scheme is capable of capturing the optimal policy $\pi^\star$ in presence of model uncertainties and disturbances.

Importantly, the deterministic policy $\pi_\theta(s)$ can be obtained as

$$\pi_\theta(s) = u_0^\star(s, \theta), \tag{19}$$

where $u_0^\star(s, \theta)$ is the first element of $u^\star$, which is the input solution of the MPC scheme (16).

### B. LSTD-Based DPG Method

The DPG method optimizes the policy parameters $\theta$ directly via gradient descent steps on the performance function $J$, defined in (11). The update rule is as follows

$$\theta \leftarrow \theta - \alpha \nabla_\theta J(\pi_\theta), \tag{20}$$

where $\alpha > 0$ is the step size. Applying the DPG method developed by [17], the gradient of $J$ with respect to parameters $\theta$ is obtained as

$$\nabla_\theta J(\pi_\theta) = \mathbb{E}\left[\nabla_\theta \pi_\theta(s) \nabla_a Q_{\pi_\theta}(s, a)|_{a=\pi_\theta}\right], \tag{21}$$

where $Q_{\pi_\theta}$ and its inner function $V_{\pi_\theta}$ are the action-value function and value function associated to the policy $\pi_\theta$, respectively, defined as follows

$$Q_{\pi_\theta}(s, a) = L(s, a) + \gamma \mathbb{E}\left[V_{\pi_\theta}(s^+|(s, a))\right] \tag{22a}$$
$$V_{\pi_\theta}(s) = Q_{\pi_\theta}(s, \pi_\theta(s)), \tag{22b}$$

where $s^+$ is the subsequent state of the state-input pair $(s, a)$. The calculations of $\nabla_\theta \pi_\theta(s)$ and $\nabla_a Q_{\pi_\theta}(s, a)$ in (21) are discussed in the following.

*1) $\nabla_\theta \pi_\theta(s)$:* The primal-dual Karush Kuhn Tucker (KKT) conditions underlying the MPC scheme (16) is written as

$$R = \begin{bmatrix} \nabla_\zeta \mathcal{L}_\theta & G_\theta & \text{diag}(\mu) H_\theta \end{bmatrix}^\top, \tag{23}$$

where $\zeta = \{\hat{\eta}, \hat{\nu}, \hat{f}, \hat{\alpha}, \sigma\}$ is the primal decision variable of the MPC (16). Term $\mathcal{L}_\theta$ is the associated Lagrange function, written as

$$\mathcal{L}_\theta(y) = \Omega_\theta + \lambda^\top G_\theta + \mu^\top H_\theta, \tag{24}$$

where $\Omega_\theta$ is the MPC cost (16a), $G_\theta$ gathers the equality constraints and $H_\theta$ collects the inequality constraints of the MPC (16). Vectors $\lambda, \mu$ are the associated dual variables. Argument $y$ reads as $y = \{\zeta, \lambda, \mu\}$ and $y^\star$ refers to the solution of the MPC (16). Consequently, the policy sensitivity $\nabla_\theta \pi_\theta$ required in (21) can then be obtained as follows ([7])

$$\nabla_\theta \pi_\theta(s) = -\nabla_\theta R(y^\star, s, \theta) \nabla_y R(y^\star, s, \theta)^{-1} \frac{\partial y}{\partial u_0}, \tag{25}$$

where $u_0$ is the first element of the input, expressed as

$$u_0 = \begin{bmatrix} \hat{f}_0^\top, \hat{\alpha}_0^\top \end{bmatrix}^\top. \tag{26}$$

*2) $\nabla_a Q_{\pi_\theta}(s, a)$:* Under some conditions [17], the action-value function $Q_{\pi_\theta}$ can be replaced by an approximator $Q_w$, i.e. $Q_w \approx Q_{\pi_\theta}$, without affecting the policy gradient. Such an approximation is labelled *compatible* and can, e.g., take the form

$$Q_w(s, a) = \underbrace{(a - \pi_\theta(s))^\top \nabla_\theta \pi_\theta(s)^\top}_{\Psi^\top(s, a)} w + V_v(s), \tag{27}$$

where $\Psi(s, a)$ is the state-action feature vector, $w$ is the parameters vector estimating the action-value function $Q_{\pi_\theta}$ and $V_v \approx V_{\pi_\theta}$ is the parameterized baseline function approximating the value function, it can take a linear form

$$V_v(s) = \Phi(s)^\top v, \tag{28}$$

where $\Phi(s)$, the state feature vector, is designed to constitute all monomials of the state with degrees less than or equal to 2. And $v$ is the corresponding parameters vector. Now we get

$$\nabla_a Q_{\pi_\theta}(s, a) \approx \nabla_a Q_w(s, a) = \nabla_\theta \pi_\theta(s)^\top w. \tag{29}$$

The parameters $w$ and $v$ of the action-value function approximation (27) are the solutions of the Least Squares (LS) problem

$$\min_{w, v} \mathbb{E}\left[(Q_{\pi_\theta}(s, a) - Q_w(s, a))^2\right], \tag{30}$$

which, in this work, is tackled via the LSTD method (see [18]). LSTD belongs to *batch method*, seeking to find the best fitting value function and action-value function, and it is more sample efficient than other methods. The LSTD update rules are as follows

$$v = \mathbb{E}_m\left\{\left[\sum_{k=1}^{K}\left[\Phi(s_k)(\Phi(s_k) - \gamma\Phi(s_{k+1}))^\top\right]\right]^{-1}\right.$$
$$\left.\sum_{k=1}^{K}\left[\Phi(s_k)L(s_k, a_k)\right]\right\}, \tag{31a}$$

$$w = \mathbb{E}_m\left\{\left[\sum_{k=1}^{K}\left[\Psi(s_k, a_k)\Psi(s_k, a_k)^\top\right]\right]^{-1}\right.$$
$$\left.\sum_{k=1}^{K}\left[(L(s_k, a_k) + \gamma V_v(s_{k+1}) - V_v(s_k))\Psi(s_k, a_k)\right]\right\}, \tag{31b}$$

where the summation is taken over the whole episode, which terminates at $K$ when the ship reaches the destination (i.e. $\|\eta_K - \eta_d\|_2^2 \leq d_{\text{error}}$). The values will be then averaged by taking expectation ($\mathbb{E}_m$) over $m$ episodes.

Finally, equation (20) can be rewritten as a compatible DPG

$$\theta \leftarrow \theta - \alpha \mathbb{E}_m\left\{\sum_{k=1}^{K}\left[\nabla_\theta \pi_\theta(s_k) \nabla_\theta \pi_\theta(s_k)^\top w\right]\right\}, \tag{32}$$

and the proposed MPC-LSTD-based DPG method is summarized in Algorithm 1.

## V. SIMULATION

In this section, we show the simulation results of an ASV freight mission problem using the introduced MPC-based RL method. We choose the initial parameters vector as $\theta_0 = \{\mathbf{0.55}, \mathbf{3}, \mathbf{3}, \mathbf{1e-7}, 60, 35, 0.5\}$, where the bold numbers represent constant vectors with suitable dimension. Other parameters values used in the simulation are given in Table I.

**Algorithm 1:** MPC-LSTD-based DPG method

**Input:** vessel model, objective function, initial parameters $\boldsymbol{\theta}_0$

**Output:** locally optimal policy $\boldsymbol{\pi}_{\boldsymbol{\theta}^\star}$

1 **repeat**
2    **for** *each episode* in $m$ episodes **do**
3       initialize $\boldsymbol{\eta}_0, \boldsymbol{\nu}_0$;
4       **while** $\|\boldsymbol{\eta}_k - \boldsymbol{\eta}_d\|_2^2 \leq d_{\text{error}}$ **do**
5          solve the MPC (16) and get $\boldsymbol{y}^\star$;
6          calculate and record the RL stage cost $L(\boldsymbol{s}_k, \boldsymbol{a}_k)$ according to (12) and the sensitivity $\nabla_{\boldsymbol{\theta}}\pi_{\boldsymbol{\theta}}(\boldsymbol{s}_k)$ according to (25);
7       **end**
8    **end**
9    calculate $\boldsymbol{v}$ according to (31a);
10    calculate $\boldsymbol{w}$ according to (31b);
11    update $\boldsymbol{\theta}$ according to (32);
12 **until** *convergence*;

TABLE I

PARAMETERS VALUES.

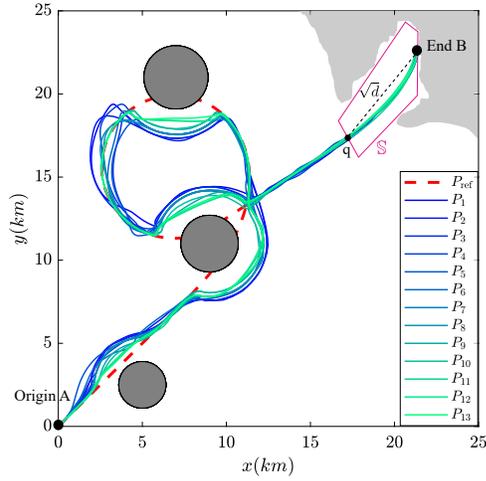| Symbol | Value | Symbol | Value |
|---|---|---|---|
| $\gamma, N, dt$ | $1, 60, 0.5$ | $\tau_a, \alpha_{\max}$ | $\mathcal{N}(0, 1e{-}3), \frac{17\pi}{18}$ |
| $f_{1\,\min,\max}$ | $-100, 100$ | $f_{2,3\,\min,\max}$ | $0, 200$ |
| $\rho, \varepsilon, \delta$ | $1, 0.001, 0.001$ | $\boldsymbol{W}$ | $\text{diag}([1,1,1])$ |
| $\boldsymbol{\omega}, \boldsymbol{\omega}_f$ | $[1, 5, 5]^\top$ | $c_{1,2,3}$ | $5, 8, 8$ |
| $d, d_s, d_{\text{error}}$ | $42.5, 1, 0.5$ | $N_o, m$ | $3, 10$ |
| $r_0, r_1, r_2, r_3$ | $1, 1.4, 1.7, 1.9$ | $\boldsymbol{\eta}_d$ | $[21.3, 23.3, 8.4]^\top$ |
| $\boldsymbol{\eta}_0$ | $[0, 0, \frac{\pi}{4}]^\top$ | $\boldsymbol{\nu}_0$ | $[0.4, 0, 0]^\top$ |



Fig. 2. Freight shipping paths from A to B. $P_{\text{ref}}$: the reference path. $P_1$-$P_{13}$: the renewed path after each learning step.
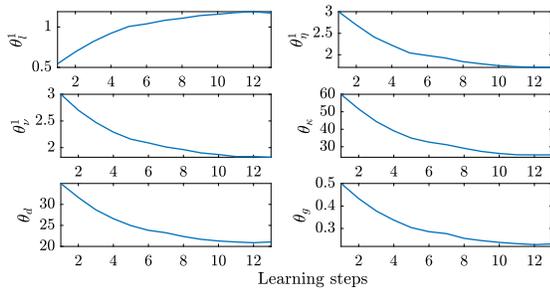


Fig. 3. Variations of some selected MPC parameters $\{\theta_l^1, \theta_\eta^1, \theta_\nu^1, \theta_\kappa, \theta_d, \theta_g\}$ over learning steps.
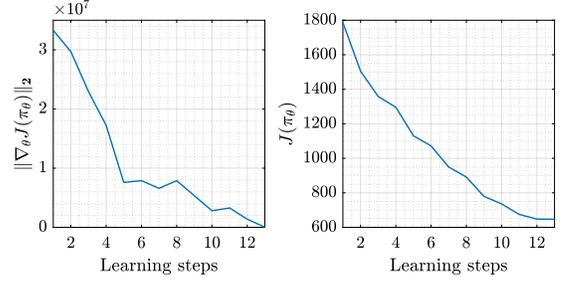


Fig. 4. Variations of the normed policy gradient $\|\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\pi}_{\boldsymbol{\theta}})\|_2$ and the closed-loop performance $J(\boldsymbol{\pi}_{\boldsymbol{\theta}})$ over learning steps.



Fig. 5. Variations of the error $\boldsymbol{\eta} - \boldsymbol{\eta}_d$ with time under the learned policy $\boldsymbol{\pi}_{\boldsymbol{\theta}^\star}$. Red line: the desired value.



Fig. 6. Variations of the vessel velocity $\boldsymbol{\nu}$ with time under the learned policy $\boldsymbol{\pi}_{\boldsymbol{\theta}^\star}$. Red line: the desired value.



Fig. 7. Variations of the thruster force $\boldsymbol{f}$ and thruster angle $\boldsymbol{\alpha}$ with time under the learned policy $\boldsymbol{\pi}_{\boldsymbol{\theta}^\star}$. Green line: the constraint value.

Figure 2 shows the prescribed reference path and the thirteen shipping paths updated after each learning step. The last path $P_{13}$ is obtained under the final learned policy $\pi_{\theta^\star}$ with an episode length of $K = 550$. It is worth noting that, although we say that if the parametrization is rich enough, the MPC scheme can generate the optimal policy, this is a theoretical result. In practice, the assumption of a "rich enough" parametrization is typically not satisfied. Other practical issues can come in the way of optimality such as, e.g., the local convergence of the RL algorithm and of the solver treating the MPC scheme. Addressing these potential issues typically requires good initial guesses. Although these are often available in the MPC context, we can only claim that the final learned policy $\pi_{\theta^\star}$ obtained from the converged parameters $\theta^\star$ is locally optimal. This observation applies to most RL techniques. Following the reference path $P_{\mathrm{ref}}$ defined from the origin $A$ to the point $q$, the vessel departs from $A$ and passes through three obstacles to reach $q$. At the point $q$, where $\|\eta_k - \eta_d\|_2^2 = d$, the vessel transits from path following to docking. The vessel eventually stops at the end $B$ with zero velocities and thruster forces, and has no collision with the quay (within the safety operation region $\mathbb{S}$) during the docking process. It can be seen that in the first few paths ($P_1$-$P_4$), the ship does not follow $P_{\mathrm{ref}}$ precisely, and is relatively far away from the three obstacles when it bypassed them. After learning, such as in the $P_{13}$, the ship follows closely the reference route, and the distance when avoiding obstacles is also reduced.

Figure 3 shows the convergences of the MPC parameters $\theta$ over learning steps ($\theta^\star$ represents the converged parameters). Note that $\theta_l^1$ is the first element of $\theta_l$, and the same fashion for others. It can be seen that the initial value of $\theta_l$ is relatively small, and the initial values of $\theta_\eta, \theta_\nu, \theta_\kappa, \theta_d$ are relatively large. Therefore, in the MPC cost (16a), the terminal cost weights more than the stage cost, i.e., docking is regarded as more important than path following. Consequently, the path following performance is relatively poor in the initial episodes, and then gets improved as $\theta_l$ increases and $\theta_\eta, \theta_\nu, \theta_\kappa, \theta_d$ decrease. In addition, the initial value of $\theta_g$ is large, which means that the ship must be very far away from the obstacles. However, this is unnecessary under the premise of ensuring the safe distance $d_s$. To reduce the cost, RL gradually reduces $\theta_g$, and therefore results in what we have in Fig. 2: the distance for avoiding obstacles tends to decrease over learning.

The variations of the normed policy gradient $\|\nabla_\theta J(\pi_\theta)\|_2$ and the closed-loop performance $J(\pi_\theta)$ are displayed in Fig. 4. As can be seen, the policy gradient converges to near zero and the performance is improved significantly over learning. Figure 5 illustrates the variations of error between the vessel pose state $\eta$ and the desired docking state $\eta_d$ under the learned policy $\pi_{\theta^\star}$. Figure 6 presents the variations of the vessel velocity $\nu$ with time under the policy $\pi_{\theta^\star}$. The red dash lines in these two figures represent the zero-valued reference lines. It can be seen that both the pose error and velocity converge to the red dash lines, which signifies a satisfactory docking. The variations of the vessel's thruster force $f$ and thruster angle $\alpha$ under policy $\pi_{\theta^\star}$ are exhibited in Fig. 7. The green lines stand for the constraint values. As can be seen, both the forces

and the angles obey their constraints, and when approaching the endpoint, the forces decline to zero and the angles remain constant.

## VI. CONCLUSION

This paper presents an MPC-based RL method for the ASV to accomplish a freight mission, which includes collision-free path following, autonomous docking, and an ingenious transition between them. We use a parameterized MPC-scheme as the policy approximation function, and adopt the LSTD-based DPG method to update the parameters such that the closed-loop performance gets improved with learning. For future works, we will further validate our proposed method by realizing the experimental implementations.

## REFERENCES

[1] J. E. Manley, "Unmanned surface vehicles, 15 years of development," in *OCEANS 2008*. IEEE, 2008, pp. 1–4.

[2] N. Gu, Z. Peng, D. Wang, Y. Shi, and T. Wang, "Antidisturbance coordinated path following control of robotic autonomous surface vehicles: Theory and experiment," *IEEE/ASME transactions on mechatronics*, vol. 24, no. 5, pp. 2386–2396, 2019.

[3] A. B. Martinsen, G. Bitar, A. M. Lekkas, and S. Gros, "Optimization-based automatic docking and berthing of asvs using exteroceptive sensors: Theory and experiments," *IEEE Access*, vol. 8, pp. 204 974–204 986, 2020.

[4] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Advances in neural information processing systems*, 2000, pp. 1057–1063.

[5] L. Buşoniu, T. de Bruin, D. Tolić, J. Kober, and I. Palunko, "Reinforcement learning for control: Performance, stability, and deep approximators," *Annual Reviews in Control*, vol. 46, pp. 8–28, 2018.

[6] M. Zanon and S. Gros, "Safe reinforcement learning using robust MPC," *IEEE Transactions on Automatic Control*, 2020.

[7] S. Gros and M. Zanon, "Data-driven economic NMPC using reinforcement learning," *IEEE Transactions on Automatic Control*, vol. 65, no. 2, pp. 636–648, 2019.

[8] E. F. Camacho and C. B. Alba, *Model predictive control*. Springer science & business media, 2013.

[9] S. Gros and M. Zanon, "Reinforcement learning for mixed-integer problems based on MPC," *ArXiv Preprint:2004.01430*, 2020.

[10] M. Zanon, S. Gros, and A. Bemporad, "Practical reinforcement learning of stabilizing economic MPC," in *2019 18th European Control Conference (ECC)*. IEEE, 2019, pp. 2258–2263.

[11] A. B. Kordabad, W. Cai, and S. Gros, "Multi-agent battery storage management using mpc-based reinforcement learning," *arXiv preprint arXiv:2106.03541*, 2021.

[12] W. Cai, H. N. Esfahani, A. B. Kordabad, and S. Gros, "Optimal management of the peak power penalty for smart grids using mpc-based reinforcement learning," *arXiv preprint arXiv:2108.01459*, 2021.

[13] A. B. Kordabad, W. Cai, and S. Gros, "MPC-based reinforcement learning for economic problems with application to battery storage," *arXiv preprint arXiv:2104.02411*, 2021.

[14] R. Skjetne, Ø. Smogeli, and T. I. Fossen, "Modeling, identification, and adaptive maneuvering of cybership ii: A complete design with experiments," *IFAC Proceedings Volumes*, vol. 37, no. 10, pp. 203–208, 2004.

[15] A. B. Martinsen, A. M. Lekkas, and S. Gros, "Autonomous docking using direct optimal control," *IFAC-PapersOnLine*, vol. 52, no. 21, pp. 97–102, 2019.

[16] T. A. Johansen, T. I. Fossen, and S. P. Berge, "Constrained nonlinear control allocation with singularity avoidance using sequential quadratic programming," *IEEE Transactions on Control Systems Technology*, vol. 12, no. 1, pp. 211–216, 2004.

[17] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *Proceedings of the 31st International Conference on International Conference on Machine Learning*. JMLR.org, 2014, p. I–387–I–395.

[18] M. G. Lagoudakis and R. Parr, "Least-squares policy iteration," *Journal of machine learning research*, vol. 4, pp. 1107–1149, 2003.