

Fault Tolerant Control for Autonomous Surface Vehicles via Model Reference Reinforcement Learning

Zhang, Qingrui; Zhang, Xinyu; Zhu, Bo ; Reppa, V.

DOI

[10.1109/CDC45484.2021.9683461](https://doi.org/10.1109/CDC45484.2021.9683461)

Publication date

2021

Document Version

Final published version

Published in

Proceedings of the 60th IEEE Conference on Decision and Control (CDC 2021)

Citation (APA)

Zhang, Q., Zhang, X., Zhu, B., & Reppa, V. (2021). Fault Tolerant Control for Autonomous Surface Vehicles via Model Reference Reinforcement Learning. In *Proceedings of the 60th IEEE Conference on Decision and Control (CDC 2021)* (pp. 1536-1541). IEEE. <https://doi.org/10.1109/CDC45484.2021.9683461>

Important note

To cite this publication, please use the final published version (if applicable). Please check the document version above.

Copyright

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

Green Open Access added to TU Delft Institutional Repository

'You share, we take care!' - Taverne project

<https://www.openaccess.nl/en/you-share-we-take-care>

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.

Fault Tolerant Control for Autonomous Surface Vehicles via Model Reference Reinforcement Learning

Qingrui Zhang¹, Xinyu Zhang¹, Bo Zhu¹, and Vasso Reppa²

Abstract—A novel fault tolerant control algorithm is proposed in this paper based on model reference reinforcement learning for autonomous surface vehicles subject to sensor faults and model uncertainties. The proposed control scheme is a combination of a model-based control approach and a data-driven method, so it can leverage the advantages of both sides. The proposed design contains a baseline controller that ensures stable tracking performance at healthy conditions, a fault observer that estimates sensor faults, and a reinforcement learning module that learns to accommodate sensor faults using fault estimation and compensate for model uncertainties. The impact of sensor faults and model uncertainties can be effectively mitigated by this composite design. Stable tracking performance can also be ensured even at both the offline training and online implementation stages for the learning-based fault tolerant control. A numerical simulation with gyro sensor faults is presented to demonstrate the efficiency of the proposed algorithm.

I. INTRODUCTION

With the impressive advancement in guidance, navigation, and control technologies, autonomous surface vehicles (ASVs) are becoming a possible alternative solution to human-operated vessels in diverse applications [1], [2]. In particular, the sudden burst of the COVID-19 pandemic makes it more urgent to develop ASVs for global shipping. In most applications, ASVs are expected to be running safely with little human intervention for a long period of time. It requires ASVs to have enough safety and reliability attributes for both avoiding catastrophic consequences and securing the deliverance of correct service [3]. However, ASVs are susceptible to malfunctions, degradation in system components, and sensor faults, *etc.*, thereby experiencing performance deterioration, instability, and even disastrous loss. Those issues motivate the extensive study of fault tolerant control (FTC), an efficient technique that can recover system performance or keep systems operational after encountering faults, and thus enhance the system safety significantly [4].

FTC algorithms are generally divided into two categories, namely passive and active FTC [5]. In passive FTC, a reliable controller is developed, which has sufficient robustness against all expected faults of low magnitude either by using a robust control approach or an adaptive control method [6]. Passive FTC algorithms demand no controller reconfiguration, so they need to accommodate both the healthy and faulty conditions

[7]. However, faults would not occur most time in real-life applications. Hence, passive FTC is conservative and has limited fault tolerant capabilities [5], [8]. Different from passive FTC, active FTC algorithms can react actively to system faults by monitoring system health using a fault diagnosis and identification (FDI) mechanism [9]. Once a fault is detected, an active fault tolerant controller can reconfigure itself efficiently to recover the system performance. Most FTC algorithms belong to model-based approaches. Passive FTC needs to know the "worst-case" system faults for the design of robust control [10]. Active FTC needs the degraded system model under faults for control reconfiguration [11].

To reduce the dependence on system modelling, reinforcement learning (RL) has been discussed as an option for the design of FTC [12]–[14]. The major advantage of RL is the learning of an optimal control law from data samples without using models. Such an advantage is very suitable for ASVs subject to significant model uncertainties and environmental disturbances [15]. However, it is demanding for model-free RL to ensure closed-loop stability, if no extra assumption is made on the initial choices of control laws. Many existing RL algorithms for FTC learn a robust optimal FTC that ensures system stability under the "worst-case" faults [16], although model information is not necessary.

In this paper, we present a novel FTC algorithm for autonomous surface vehicles subject to sensor faults and model uncertainties based on reinforcement learning. Via the integration of RL with a model-based control approach, the proposed control scheme, termed as model reference reinforcement learning [17], [18], can ensure the closed-loop stability. It can also avoid learning the "worst-case" controller by incorporating a fault diagnosis and estimation mechanism to signal the occurrence and magnitude of sensor faults. Hence, our FTC algorithm can actively react to sensor faults, mitigate the impact of both sensor faults and model uncertainties on the trajectory tracking performance of ASVs, and ensure stable and safe tracking performance at both the offline training and the online implementation stages. In summary, two major contributions of this work are identified.

- 1) A new formulation framework following the model-reference structure is presented for the design of reinforcement learning-based control. With this new formulation, it is possible to combine model-based approaches with data-driven methods, and leverage the advantages from both sides.
- 2) A novel RL-based active FTC algorithm is presented for ASVs subject to sensor faults and model uncertainties. Numerical simulations show that the proposed FTC

¹School of Aeronautics and Astronautics, Sun Yat-sen University, Guangzhou, Guangdong, China (zhangqr9@mail.sysu.edu.cn; zhangxy385@mail2.sysu.edu.cn; zhubo5@mail.sysu.edu.cn)

²Department of Maritime and Transport Technology, Delft University of Technology, Delft, the Netherlands (V.Reppa@tudelft.nl)

algorithm can efficiently mitigate the influence of sensor faults and model uncertainties.

The rest of the paper is organized as follows. Problem formulation is provided in Section II. Section III presents the reinforcement learning-based FTC scheme. In Section IV, we describe details on the real-life implementation of the proposed algorithm. Numerical simulations are provided in Section V. Concluding remarks are given in Section VI.

II. PROBLEM FORMULATION

According to [19], [20], the ASV dynamics are

$$\begin{cases} \dot{\boldsymbol{\eta}} = \mathcal{R}(\boldsymbol{\eta}) \boldsymbol{\nu} \\ \mathcal{M} \dot{\boldsymbol{\nu}} + (\mathcal{C}(\boldsymbol{\nu}) + \mathcal{D}(\boldsymbol{\nu})) \boldsymbol{\nu} + \mathcal{G}(\boldsymbol{\nu}) = \mathcal{B} \mathbf{u} \end{cases} \quad (1)$$

where $\boldsymbol{\eta} = [x_p, y_p, \psi_p]^T \in \mathbb{R}^3$ is a generalized coordinate vector with x_p and y_p denoting the horizontal position coordinates of an ASV in the inertial frame and ψ_p the heading angle, $\boldsymbol{\nu} = [u_p, v_p, r_p]^T \in \mathbb{R}^3$ is the generalized speed vector with u_p and v_p being the linear velocities in surge (x -axis) and sway (y -axis), respectively, and r_p the heading angular rate, $\mathbf{u} = [\tau_u, \tau_r]^T \in \mathbb{R}^2$ is the control forces and moments, $\mathcal{G}(\boldsymbol{\nu}) = [\mathbf{g}_1(\boldsymbol{\nu}), \mathbf{g}_2(\boldsymbol{\nu}), \mathbf{g}_3(\boldsymbol{\nu})]^T \in \mathbb{R}^3$ is unmodeled dynamics due to gravitational and buoyancy forces and moments [19], $\mathcal{B} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}^T$ is the input matrix. $\mathcal{M} = \mathcal{M}^T \in \mathbb{R}^{3 \times 3}$ is the inertia matrix,

$$\mathcal{M} = [\mathcal{M}_{ij}] = \begin{bmatrix} \mathcal{M}_{11} & 0 & 0 \\ 0 & \mathcal{M}_{22} & \mathcal{M}_{23} \\ 0 & \mathcal{M}_{32} & \mathcal{M}_{33} \end{bmatrix} \quad (2)$$

where $\mathcal{M}_{11} = m - \mathcal{X}_{\dot{u}}$, $\mathcal{M}_{22} = m - \mathcal{Y}_{\dot{v}}$, $\mathcal{M}_{33} = I_z - \mathcal{N}_{\dot{r}}$, and $\mathcal{M}_{32} = \mathcal{M}_{23} = m x_g - \mathcal{Y}_{\dot{r}}$. The matrix $\mathcal{C}(\boldsymbol{\nu}) = -\mathcal{C}^T(\boldsymbol{\nu})$ contains the Coriolis and centripetal terms, so

$$\mathcal{C} = [\mathcal{C}_{ij}] = \begin{bmatrix} 0 & 0 & \mathcal{C}_{13}(\boldsymbol{\nu}) \\ 0 & 0 & \mathcal{C}_{23}(\boldsymbol{\nu}) \\ -\mathcal{C}_{13}(\boldsymbol{\nu}) & -\mathcal{C}_{23}(\boldsymbol{\nu}) & 0 \end{bmatrix} \quad (3)$$

where $\mathcal{C}_{13}(\boldsymbol{\nu}) = -\mathcal{M}_{22}v - \mathcal{M}_{23}r$, $\mathcal{C}_{23}(\boldsymbol{\nu}) = \mathcal{M}_{11}u$. The damping matrix $\mathcal{D}(\boldsymbol{\nu})$ is

$$\mathcal{D}(\boldsymbol{\nu}) = [\mathcal{D}_{ij}] = \begin{bmatrix} \mathcal{D}_{11}(\boldsymbol{\nu}) & 0 & 0 \\ 0 & \mathcal{D}_{22}(\boldsymbol{\nu}) & \mathcal{D}_{23}(\boldsymbol{\nu}) \\ 0 & \mathcal{D}_{32}(\boldsymbol{\nu}) & \mathcal{D}_{33}(\boldsymbol{\nu}) \end{bmatrix} \quad (4)$$

where $\mathcal{D}_{11}(\boldsymbol{\nu}) = -\mathcal{X}_u - \mathcal{X}_{|u|}|u| - \mathcal{X}_{uuu}u^2$, $\mathcal{D}_{22}(\boldsymbol{\nu}) = -\mathcal{Y}_v - \mathcal{Y}_{|v|}|v| - \mathcal{Y}_{|r|v}|r|$, $\mathcal{D}_{23}(\boldsymbol{\nu}) = -\mathcal{Y}_r - \mathcal{Y}_{|v|r}|v| - \mathcal{Y}_{|r|r}|r|$, $\mathcal{D}_{32}(\boldsymbol{\nu}) = -\mathcal{N}_v - \mathcal{N}_{|v|}|v| - \mathcal{N}_{|r|v}|r|$, $\mathcal{D}_{33}(\boldsymbol{\nu}) = -\mathcal{N}_r - \mathcal{N}_{|v|r}|v| - \mathcal{N}_{|r|r}|r|$, and $\mathcal{X}(\cdot)$, $\mathcal{Y}(\cdot)$, and $\mathcal{N}(\cdot)$ are hydrodynamic coefficients [19]. The rotation matrix \mathcal{R} is

$$\mathcal{R}(\boldsymbol{\eta}) = \begin{bmatrix} \cos \psi_p & -\sin \psi_p & 0 \\ \sin \psi_p & \cos \psi_p & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

By defining $\mathbf{x} = [\boldsymbol{\eta}^T, \boldsymbol{\nu}^T]^T$, it yields

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & \mathcal{R}(\boldsymbol{\eta}) \\ 0 & \mathcal{H}(\boldsymbol{\nu}) \end{bmatrix} \mathbf{x} + \begin{bmatrix} 0 \\ \mathcal{N} \end{bmatrix} \mathbf{u} + \begin{bmatrix} 0 \\ -\mathcal{M}^{-1} \mathcal{G}(\boldsymbol{\nu}) \end{bmatrix} \quad (5)$$

where $\mathcal{H}(\boldsymbol{\nu}) = -\mathcal{M}^{-1}(\mathcal{C}(\boldsymbol{\nu}) + \mathcal{D}(\boldsymbol{\nu}))$ and $\mathcal{N} = -\mathcal{M}^{-1} \mathcal{B}$. The state measurement of the ASV system (1) is corrupted by noises and sensor faults, so it is expressed as

$$\mathbf{y} = \mathbf{x} + \mathbf{n} + \mathbf{f}(t)$$

where $\mathbf{n} \in \mathbb{R}^6$ is the measurement noise and $\mathbf{f}(t) \in \mathbb{R}^6$ denotes the possible sensor fault. In this paper, we only consider sensor faults on the measurement of the heading angular rate r_p , so $\mathbf{f}(t) = [0, 0, 0, 0, 0, f_r(t)]^T$. The sensor fault $f_r(t)$ is given by

$$f_r(t) = \beta(t - T_f) \phi(t - T_f)$$

where $\phi(t - T_f)$ is the unknown function of the sensor fault that occurs at the time instant T_f , and $\beta(t - T_f)$ is the time profile with $\beta(t - T_f) = 0$ for $t \leq T_f$ and $\beta(t - T_f) = 1 - e^{-k(t - T_f)}$ for $t > T_f$ where k is the evolution rate of the fault [21], [22]. Note that $k \rightarrow \infty$, if the occurrence of a sensor fault is abrupt, e.g., bias fault.

III. REINFORCEMENT LEARNING-BASED FAULT TOLERANT CONTROL SCHEME

This section starts with the presentation of a model-reference control structure, then provides preliminaries on RL, and eventually, presents the new FTC scheme.

A. Model-reference control structure

For most ASV systems, accurate nonlinear dynamic model is rarely available. Major uncertainties come from \mathcal{M} , $\mathcal{C}(\boldsymbol{\nu})$, and $\mathcal{D}(\boldsymbol{\nu})$ due to hydrodynamics, and $\mathcal{G}(\boldsymbol{\nu})$ due to gravitational and buoyancy forces and moments. Although the ASV dynamics are subject to uncertainties, a nominal model is still available based on the known information on the ASV dynamics (5). The nominal model of (5) is

$$\dot{\mathbf{x}}_m = \begin{bmatrix} 0 & \mathcal{R}(\boldsymbol{\eta}_m) \\ 0 & \mathcal{H}_m \end{bmatrix} \mathbf{x}_m + \begin{bmatrix} 0 \\ \mathcal{N}_m \end{bmatrix} \mathbf{u}_m \quad (6)$$

where \mathcal{N}_m and \mathcal{H}_m contain all the known constant parameters of the ASV dynamics (5), and $\boldsymbol{\eta}_m = [x_m, y_m, \psi_m]^T \in \mathbb{R}^3$ is a generalized coordinate vector of the nominal model. In this paper, \mathcal{M}_m is given by $\mathcal{M}_m = \text{diag}\{\mathcal{M}_{11}, \mathcal{M}_{22}, \mathcal{M}_{33}\}$, $\mathcal{H}_m = \mathcal{M}_m^{-1} \mathcal{D}_m$ with $\mathcal{D}_m = \text{diag}\{-\mathcal{X}_u, -\mathcal{Y}_v, -\mathcal{N}_r\}$, and $\mathcal{N}_m = \mathcal{M}_m^{-1} \mathcal{B}$. Hence, in the nominal model, all the nonlinear terms in the inner-loop dynamics are ignored, so we end up with a linear and decoupled model for the dynamic equations of the generalized velocity vector $\boldsymbol{\nu}$. As the dynamics of the nominal model (6) are known, it is possible to design a control law \mathbf{u}_m allowing the states of the nominal system (6) to converge to a reference signal \mathbf{x}_r , i.e., $\|\mathbf{x}_m - \mathbf{x}_r\|_2 \rightarrow 0$ as $t \rightarrow \infty$. The control law \mathbf{u}_m can also be used by the full ASV dynamics (5) as a baseline control.

In the model-reference control structure, the objective is to design a control law allowing the states of (5) to track those of the nominal model (6), so the overall control law is

$$\mathbf{u} = \mathbf{u}_b + \mathbf{u}_l \quad (7)$$

where \mathbf{u}_b is a baseline controller, and \mathbf{u}_l is a control policy from the deep RL module. The baseline control \mathbf{u}_b is employed to ensure some basic performance, i.e., local stability, while \mathbf{u}_l is employed to compensate for system uncertainties and sensor faults. The baseline control \mathbf{u}_b is designed based on the nominal model (6), so it has the same expression as \mathbf{u}_m in (6). Note that \mathbf{u}_b or \mathbf{u}_m can be designed by any existing model-based method. Hence, we focus on the development of \mathbf{u}_l using RL.

B. Reinforcement learning

The formulation of RL is based on a Markov decision process (MDP) denoted by a tuple $MDP := \langle \mathcal{S}, \mathcal{U}, \mathcal{P}, R, \gamma \rangle$, where \mathcal{S} is the state space, \mathcal{U} specifies the action/input space, $\mathcal{P} : \mathcal{S} \times \mathcal{U} \times \mathcal{S} \rightarrow \mathbb{R}$ defines a transition probability, $R : \mathcal{S} \times \mathcal{U} \rightarrow \mathbb{R}$ is a reward function, and $\gamma \in [0, 1)$ is a discount factor. In MDP, $s \in \mathcal{S}$ contains all available signals affecting the RL control $u_l \in \mathcal{U}$. In this paper, the transition probability is characterized by (1) and a reference signal, x_r .

Let s_t be the state signal s at the time step t , and accordingly, $u_{l,t}$ be the input from the RL-based control. The RL algorithm aims to maximize an action-value function, a.k.a., Q-function, given as

$$Q_\pi(s_t, u_{l,t}) = R_t + \gamma \mathbb{E}_{s_{t+1}} [V_\pi(s_{t+1})] \quad (8)$$

where R_t is the reward function with $R_t = R(s_t, u_{l,t})$, $\mathbb{E}_{s_{t+1}} [V_\pi(s_{t+1})] = \sum_{s_{t+1}} \mathcal{P}_{t+1|t} [V_\pi(s_{t+1})]$, and $V_\pi(s_{t+1})$ is called state value function for s_{t+1} , where

$$\begin{aligned} V_\pi(s_t) &= \sum_{u_{l,t}} \pi(u_{l,t}|s_t) \mathbb{E}_{s_{t+1}} [R_t + \gamma V_\pi(s_{t+1}) + \\ &\quad + \alpha \mathcal{H}(\pi(u_{l,t}|s_t))] \quad (9) \\ &= \mathbb{E}_\pi [\mathbb{E}_{s_{t+1}} [R_t + \gamma V_\pi(s_{t+1})] - \alpha \log \pi(u_{l,t}|s_t)] \end{aligned}$$

where $\mathcal{H}(\pi(u_{l,t}|s_t)) = -\mathbb{E}_\pi [\log(\pi(u_{l,t}|s_t))]$ is the entropy of the policy, α is a temperature parameter, and $\pi(u_{l,t}|s_t)$ is the control policy that is the probability of choosing an action $u_{l,t} \in \mathcal{U}$ at a state $s_t \in \mathcal{S}$ [18].

The objective in RL is to solve the optimization problem below.

$$\pi^* = \arg \max Q_\pi(s_t, u_{l,t}) \quad (10)$$

C. Fault diagnosis and estimation

For the clarity of the presentation, we only consider sensor faults in the measurement of inner-loop states that are ν . Hence, only the inner-loop dynamics of ASVs are considered in the development of the fault diagnosis and estimation. However, the proposed algorithm can be extended to more generic situations, for instance, sensor faults in the measurement of positions.

The overall control design is based on the model-reference control structure given in Section III-A, so the uncertain inner-loop dynamics of the ASV model (5) is rewritten as

$$\dot{\nu} = \mathcal{H}_m \nu + \mathcal{N}_m (u_b + u_l) + \beta(\nu) \quad (11)$$

where $\beta(\nu)$ is the aggregation of all uncertainties in the inner-loop dynamics. Assume that $\beta(\nu)$ is bounded. Let $e_\nu = \nu - \nu_m$. According to (6) and (12), one has

$$\dot{e}_\nu = \mathcal{H}_m e_\nu + \mathcal{N}_m (u_b - u_m) + \mathcal{N}_m u_l + \beta(\nu) \quad (12)$$

Under healthy conditions, the model uncertainty term $\beta(\nu)$ can be fully compensated using a learning-based control u_l according to [17], [18]. It implies that $\|e_\nu(t)\|_2 \leq \epsilon$ as $t \rightarrow \infty$, where ϵ is a certain positive small constant. If sensor faults happened, the error signal e_ν will be large than ϵ . A naive idea for the learning-based FTC is to treat the sensor faults

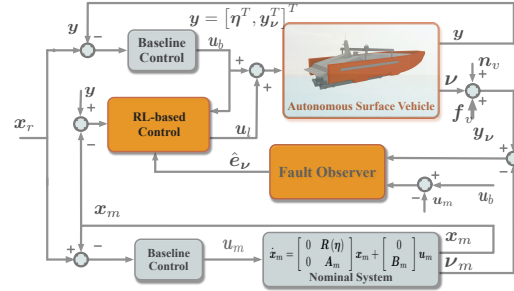


Fig. 1: RL-based FTC scheme

as part of external disturbances. However, treating sensor faults as disturbances will result in a conservative learning-based control like the robust control. Hence, we introduce a fault diagnosis and estimation mechanism that allows the learning-based control to adapt to different scenarios: healthy and unhealthy conditions.

Let $y_\nu = \nu + n_\nu + f_\nu$, where n_ν denotes the measurement noise and f_ν is the sensor fault. Furthermore, we define $e_{y_\nu} = y_\nu - \nu_m = e_\nu + n_\nu + f_\nu$ as the faulty residual vector. In real applications, e_{y_ν} is measurable instead of e_ν . The fault diagnosis and estimation mechanism is, therefore,

$$\dot{\hat{e}_\nu} = \mathcal{H}_m \hat{e}_\nu + \mathcal{N}_m (u_b - u_m) + \mathbf{L} (e_{y_\nu} - \hat{e}_\nu) \quad (13)$$

where \mathbf{L} is chosen such that $\mathcal{H}_m - \mathbf{L}$ is Hurwitz. The signal \hat{e}_ν performs as the indicator of the occurrence and strength of the sensor faults. Let $\varepsilon_\nu = e_\nu - \hat{e}_\nu$, and we have

$$\dot{\varepsilon}_\nu = (\mathcal{H}_m - \mathbf{L}) \varepsilon_\nu + \mathcal{N}_m u_l + \beta(\nu) + \mathbf{L} (n_\nu + f_\nu) \quad (14)$$

The following lemma exists for (14)

Lemma 1: Suppose that the model uncertainty term $\beta(\nu)$ can be fully compensated using a RL-based control u_l . The error dynamic model (14) is input-to-state stable with respect to the sensor fault f_ν .

Proof: Let $\varepsilon_\nu = \varepsilon_\nu^1 + \varepsilon_\nu^2$ with

$$\begin{aligned} \dot{\varepsilon}_\nu^1 &= (\mathcal{H}_m - \mathbf{L}) \varepsilon_\nu^1 + \mathcal{N}_m u_l + \beta(\nu) \\ \dot{\varepsilon}_\nu^2 &= (\mathcal{H}_m - \mathbf{L}) \varepsilon_\nu^2 + \mathbf{L} (n_\nu + f_\nu) \end{aligned}$$

If $\beta(\nu)$ is fully compensated by u_l , it implies that $\mathcal{N}_m u_l + \beta(\nu)$ together can be treated as a bounded negligible external signal. Since $\mathcal{H}_m - \mathbf{L}$ is Hurwitz, ε_ν^1 will be input-to-state stable with respect to $\mathcal{N}_m u_l + \beta(\nu)$. It implies that ε_ν^1 will be negligible as well. Hence, the magnitude of ε_ν mainly results from ε_ν^2 . Since $\mathcal{H}_m - \mathbf{L}$ is Hurwitz, ε_ν^2 is input-to-state stable with respect to f_ν . Hence, we can conclude that ε_ν is input-to-state stable with respect to f_ν . ■

D. RL-based fault tolerant control

The RL-based fault tolerant control is developed using the output from the fault diagnosis and estimation mechanism (13) as shown in Figure 1. RL learns the control policies using data samples, including input and state data, at discrete time steps. Assume that the sample time step is fixed and denoted by δt . Without loss of generality, let y_t , $u_{b,t}$, $u_{l,t}$, and $\hat{e}_{\nu,t}$ be the ASV state, the baseline control action, the control action

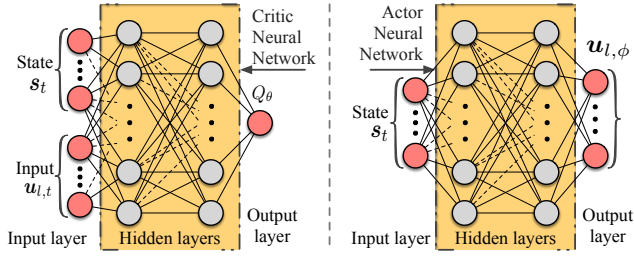


Fig. 2: MLP networks for Q_θ and π_ϕ

from RL, and the output of the fault diagnosis and estimation mechanism at the time step t , respectively. The state signal s at the time step t is thus $s_t = \begin{bmatrix} \mathbf{x}_{m,t}^T - \mathbf{y}_t^T, \mathbf{u}_{b,t}^T, \hat{\mathbf{e}}_{\nu,t}^T \end{bmatrix}^T$.

The learning-based fault tolerant robust control \mathbf{u}_l will be learned instead of designed based on RL according to Section III-B. The training/learning process of RL will repeatedly execute policy evaluation and policy improvement. In the policy evaluation, the Q-value in (8) is computed by applying a Bellman operation $Q_\pi(s_t, \mathbf{u}_{l,t}) = \mathcal{T}^\pi Q_\pi(s_t, \mathbf{u}_{l,t})$ where

$$\mathcal{T}^\pi Q_\pi(s_t, \mathbf{u}_{l,t}) = R_t + \gamma \mathbb{E}_{s_{t+1}} \left\{ \mathbb{E}_\pi [Q_\pi(s_{t+1}, \mathbf{u}_{l,t+1}) - \alpha \ln(\pi(\mathbf{u}_{l,t+1}|s_{t+1}))] \right\} \quad (15)$$

In the policy improvement, the policy is updated by

$$\pi_{new} = \arg \min_{\pi' \in \Pi} \mathcal{D}_{KL} \left(\pi'(\cdot|s_t) \parallel Z^{\pi_{old}} e^{\frac{1}{\alpha} Q^{\pi_{old}}(s_t, \cdot)} \right) \quad (16)$$

where Π denotes a policy set, π_{old} denotes the policy from the last update, $Q^{\pi_{old}}$ is the Q-value of π_{old} , \mathcal{D}_{KL} denotes the Kullback-Leibler (KL) divergence, and $Z^{\pi_{old}}$ is a normalization factor. The objective can be transformed into

$$\pi^* = \arg \min_{\pi \in \Pi} \mathbb{E}_\pi \left[\alpha \ln(\pi(\mathbf{u}_{l,t}|s_t)) - Q(s_t, \mathbf{u}_{l,t}) \right] \quad (17)$$

More details on how (17) is obtained can be found in [23].

IV. ALGORITHM IMPLEMENTATION USING DEEP NEURAL NETWORKS

In this paper, both $Q_\pi(s_t, \mathbf{u}_{l,t})$ and $\pi(\mathbf{u}_{l,t}|s_t)$ are approximated by fully connected multiple layer perceptrons (MLP) with rectified linear unit (ReLU) nonlinearities as the activation functions. The ReLU function is

$$\text{relu}(z) = \max\{z, 0\}$$

The ReLU activation function outperforms other activation functions like sigmoid functions [24]. For a vector $z = [z_1, \dots, z_n]^T$, $\text{relu}(z) = [\text{relu}(z_1), \dots, \text{relu}(z_n)]^T$. More details on the MLP can be found in [18].

The whole training process will be offline. At each time step $t+1$, we collect data samples, such as an input from the last time step $\mathbf{u}_{l,t}$, a state from the last time step s_t , a reward R_t , and a current state s_{t+1} . Those historical data will be stored as a tuple $(s_t, \mathbf{u}_{l,t}, R_t, s_{t+1})$ at a replay memory \mathcal{D} [25]. At each policy evaluation or improvement step, we randomly sample a batch of historical data, \mathcal{B} , from the replay memory \mathcal{D} for the training of the parameters θ and ϕ . Starting the training, we apply the baseline control policy

\mathbf{u}_b to an ASV system to collect the initial data \mathcal{D}_0 as shown in Algorithm 1. The initial data set \mathcal{D}_0 is used for the initial fitting of Q-value functions. When the initialization is over, we execute both \mathbf{u}_b and the latest updated reinforcement learning policy $\pi_\phi(\mathbf{u}_{l,t}|s_t)$ to run the ASV system.

The parameters θ are trained to minimize

$$J_Q(\theta) = \mathbb{E}_{(s_t, \mathbf{u}_{l,t}) \sim \mathcal{D}} \left[\frac{1}{2} (Q_\theta(s_t, \mathbf{u}_{l,t}) - Y_{target})^2 \right] \quad (18)$$

where $(s_t, \mathbf{u}_{l,t}) \sim \mathcal{D}$ implies that we randomly pick data samples $(s_t, \mathbf{u}_{l,t})$ from a replay memory \mathcal{D} , and

$$Y_{target} = R_t + \gamma \mathbb{E}_{s_{t+1}} \left[\mathbb{E}_\pi [Q_{\bar{\theta}}(s_{t+1}, \mathbf{u}_{l,t+1}) - \alpha \log(\pi_\phi)] \right]$$

where $\bar{\theta}$ is the target parameter which will be updated slowly. The DNN parameters θ are obtained by applying the stochastic gradient descent to (18) on a data batch \mathcal{B} with a fixed size denoted by $|\mathcal{B}|$. In the final implementation, we use two critics which are parameterized by θ_1 and θ_2 , respectively. The two critics are introduced to reduce the over-estimation issue in the training of critic neural networks [26], so

$$Y_{target} = R_t + \gamma \min \left\{ Q_{\bar{\theta}_1}(s_{t+1}, \mathbf{u}_{l,t+1}), Q_{\bar{\theta}_2}(s_{t+1}, \mathbf{u}_{l,t+1}) \right\} - \gamma \alpha \log(\pi_\phi) \quad (19)$$

The policy improvement is to minimize the following objective function using data samples from the replay memory.

$$J_\pi(\phi) = \mathbb{E}_{(s_t, \mathbf{u}_{l,t}) \sim \mathcal{D}} \left(\alpha \log(\pi_\phi) - Q_\theta(s_t, \mathbf{u}_{l,t}) \right) \quad (20)$$

Parameter ϕ is trained using a stochastic gradient descent technique. At the training stage, the actor neural network is

$$\mathbf{u}_{l,\phi} = \bar{\mathbf{u}}_{l,\phi} + \sigma_\phi^{\frac{1}{2}} \odot \boldsymbol{\xi} \quad (21)$$

where $\bar{\mathbf{u}}_{l,\phi}$ is the parameterized control law to be learned, $\sigma_\phi^{\frac{1}{2}}$ is the standard deviation of the exploration noise, $\boldsymbol{\xi} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is the Gaussian noise, and “ \odot ” is the Hadamard product. The exploration noise $\boldsymbol{\xi}$ is only applied to the training stage. Once the training is done, we only need $\bar{\mathbf{u}}_{l,\phi}$ in the implementation. Hence, at the training stage, \mathbf{u}_l in Figure 1 is equal to $\mathbf{u}_{l,\phi}$. Once the training is over, we have $\mathbf{u}_l = \bar{\mathbf{u}}_{l,\phi}$.

The temperature parameter α is also updated by minimizing

$$J_\alpha = \mathbb{E}_\pi \left[-\alpha \log \pi(\mathbf{u}_{l,t}|s_t) - \alpha \bar{\mathcal{H}} \right] \quad (22)$$

where $\bar{\mathcal{H}}$ is a target entropy. The entire algorithm is summarized in Algorithm 1, where ι_Q , ι_π , and ι_α are positive learning rates (scalars), and $\kappa > 0$ is a constant scalar.

V. NUMERICAL SIMULATIONS

In this section, the proposed learning-based control algorithm is implemented to the trajectory tracking control of a supply ship model presented in [20]. Model parameters can be found in [18]. The unmodeled dynamics in the simulations are given by $g_1 = 0.279uv^2 + 0.342v^2r$, $g_2 = 0.912u^2v$, and $g_3 = 0.156ur^2 + 0.278urv^3$, respectively. The baseline control \mathbf{u}_b is designed based on a nominal model in (6) in

Algorithm 1 Reinforcement learning for fault tolerant control

- 1: Initialize parameters θ_1, θ_2 for Q_{θ_1} and Q_{θ_2} , respectively, and ϕ for the actor network (21).
- 2: Assign values to the the target parameters $\bar{\theta}_1 \leftarrow \theta_1, \bar{\theta}_2 \leftarrow \theta_2, \mathcal{D} \leftarrow \emptyset, \mathcal{D}_0 \leftarrow \emptyset$,
- 3: Get data set \mathcal{D}_0 by running \mathbf{u}_b on (5) with $\mathbf{u}_l = \mathbf{0}$
- 4: Turn off the exploration and train initial critic parameters θ_1^0, θ_2^0 using \mathcal{D}_0 according to (18).
- 5: Initialize the replay memory $\mathcal{D} \leftarrow \mathcal{D}_0$
- 6: Assign initial values to critic parameters $\theta_1 \leftarrow \theta_1^0, \theta_2 \leftarrow \theta_2^0$ and their targets $\bar{\theta}_1 \leftarrow \theta_1^0, \bar{\theta}_2 \leftarrow \theta_2^0$
- 7: **repeat**
- 8: **for** each data collection step **do**
- 9: Choose an action $\mathbf{u}_{l,t}$ according to $\pi_\phi(\mathbf{u}_{l,t}|\mathbf{s}_t)$
- 10: Run both (5), (6), and (13) & collect $\mathbf{s}_{t+1} = \{\mathbf{x}_{t+1}, \mathbf{x}_{m,t+1}, \mathbf{u}_{b,t+1}\}$
- 11: $\mathcal{D} \leftarrow \mathcal{D} \cup \{\mathbf{s}_t, \mathbf{u}_{l,t}, R(\mathbf{s}_t, \mathbf{u}_{l,t}), \mathbf{s}_{t+1}\}$
- 12: **end for**
- 13: **for** each gradient update step **do**
- 14: Sample a batch of data \mathcal{B} from \mathcal{D}
- 15: $\theta_j \leftarrow \theta_j - \iota_Q \nabla_\theta J_Q(\theta_j)$, and $j = 1, 2$
- 16: $\phi \leftarrow \phi - \iota_\pi \nabla_\phi J_\pi(\phi)$,
- 17: $\alpha \leftarrow \alpha - \iota_\alpha \nabla_\alpha J_\alpha(\alpha)$
- 18: $\bar{\theta}_j \leftarrow \kappa \bar{\theta}_j + (1 - \kappa) \theta_j$, and $j = 1, 2$
- 19: **end for**
- 20: **until** convergence (i.e. $J_Q(\theta) < \text{a small threshold}$)

TABLE I: RL configurations

Parameters	Values
Learning rate ι_Q	0.001
Learning rate ι_π	0.0001
Learning rate ι_α	0.0001
κ	0.01
Actor neural network	fully connected with three hidden layers (256 × 128 × 64 neurons)
critic neural networks	fully connected with two hidden layers (256 × 256 × 32 neurons)
Replay memory capacity	1.5×10^6
Sample batch size	512
γ	0.998
Training episodes	1500
Steps per episode	1000
time step size δt	0.1

terms of the PID control method. The reference signal is assumed to be produced by the following motion planner,

$$\dot{\boldsymbol{\eta}}_r = \mathbf{R}(\boldsymbol{\eta}_r) \boldsymbol{\nu}_r \quad \dot{\boldsymbol{\nu}}_r = \mathbf{a}_r \quad (23)$$

where $\boldsymbol{\eta}_r = [x_r, y_r, \psi_r]^T$, $\boldsymbol{\nu}_r = [u_r, 0, r_r]^T$, and $\mathbf{a}_r = [\dot{u}_r, 0, \dot{r}_r]^T$. The initial position vector is $\boldsymbol{\eta}_r(0) = [0, 0, \frac{\pi}{4}]^T$. We set $u_r(0) = 0.4 \text{ m/s}$, $r_r(0) = 0 \text{ rad/s}$, and $\dot{u}_r = 0$. The derivative of the reference angular rate \dot{r}_r is

$$\dot{r}_r = \begin{cases} \frac{\pi}{300} \text{ rad/s}^2 & \text{if } 25 \text{ s} \leq t < 35 \text{ s} \\ -\frac{\pi}{300} \text{ rad/s}^2 & \text{if } 35 \text{ s} \leq t < 45 \text{ s} \\ -\frac{\pi}{300} \text{ rad/s}^2 & \text{if } 65 \text{ s} \leq t < 75 \text{ s} \\ \frac{\pi}{300} \text{ rad/s}^2 & \text{if } 75 \text{ s} \leq t < 85 \text{ s} \\ 0 \text{ rad/s}^2 & \text{otherwise} \end{cases} \quad (24)$$

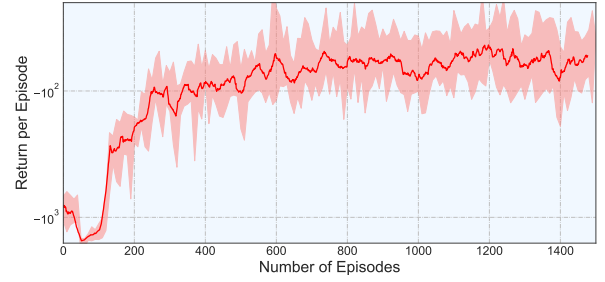


Fig. 3: Learning curves of the RL-based FTC

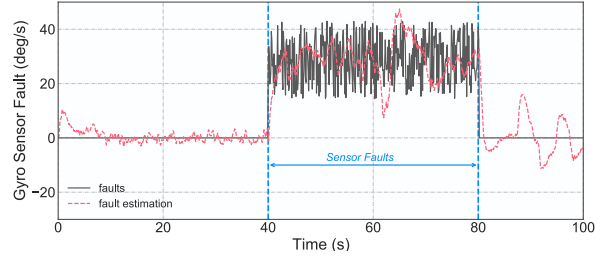


Fig. 4: Sensor fault and its estimation

The bias sensor faults of interest are defined as

$$f_r(t) = \beta(t - T_f)(\phi(t - T_f) + n_\phi) \quad (25)$$

where n_ϕ is a random noise with a uniformly random distribution. In the simulation, Gaussian measurement noises $\mathbf{n} \sim \mathcal{N}(0, \boldsymbol{\sigma}_n)$ are added to the measurement, where $\boldsymbol{\sigma}_n = \text{diag}\{0.025, 0.025, 0.01, 0.01, 0.005, 0.005\}$.

At the training stage, we run the ASV system for 100 s, and the training processes are repeated for 1500 times (i.e., 1500 episodes). Figure 3 shows the learning curves of the proposed algorithm. At each episode, we uniformly randomly sample $x(0)$ and $y(0)$ from $(-1.5, 1.5)$, $\psi(0)$ from $(-0.25\pi, 0.25\pi)$ and $u(0)$ from $(0.1, 0.6)$, and we choose $v(0) = 0$ and $r(0) = 0$. The proposed control algorithm is compared with

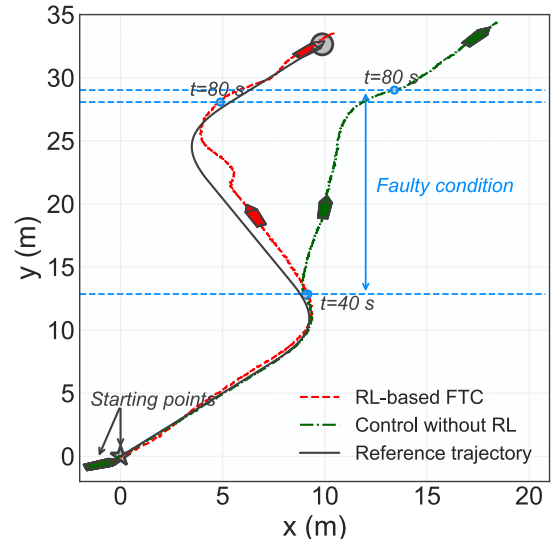


Fig. 5: Trajectory tracking performance

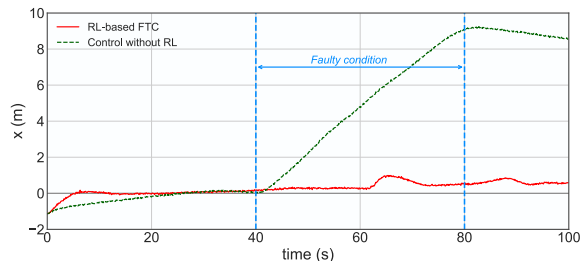


Fig. 6: Position tracking error in the x coordinate

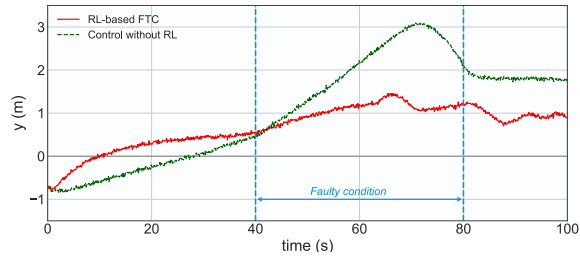


Fig. 7: Position tracking error in the y coordinate

a benchmark design in which only the baseline control u_b is considered. Configurations for the training and neural networks are found in Table I. The matrix G and H are chosen to be $G = \text{diag}\{0.025, 0.025, 0.0016, 0.005, 0.001, 0\}$ and $H = \text{diag}\{1.25e^{-4}, 1.25e^{-4}, 8.3e^{-5}\}$, respectively.

At the evaluation stage, the sensor fault profile shown in Fig. 4 is implemented to the ASV. The simulation results are summarized in Figs. 5–7. With our new design, the ASV learns to adapt to a bias sensor fault that happens to the measurement of the angular rate r_p . The trajectory tracking under faulty scenario is significantly improved.

VI. CONCLUSIONS

In this paper, a novel reinforcement learning-based fault tolerant control algorithm is presented for ASV systems subject to model uncertainties and sensor faults. The new algorithm is obtained by combining a model-reference reinforcement learning with a fault diagnosis and estimation mechanism. With the new RL-based fault tolerant control, we ensured the ASV can learn to adapt to bias faults in the gyro and recover the trajectory tracking performance under faulty conditions.

REFERENCES

- [1] P. Švec, A. Thakur, E. Raboin, B. C. Shah, and S. K. Gupta, "Target following with motion prediction for unmanned surface vehicle operating in cluttered environments," *International Journal of Robotics Research*, no. 36, pp. 383–405, Apr. 2014.
- [2] D. D. Bloisi, F. Previtali, A. Pennisi, D. Nardi, and M. Fiorini, "Enhancing automatic maritime surveillance systems with visual information," *IEEE Transactions on Intelligent Transportation Systems*, vol. 18, no. 4, pp. 824–833, Apr. 2017.
- [3] A. Avizienis, J.-C. Laprie, B. Randell, and C. Landwehr, "Basic concepts and taxonomy of dependable and secure computing," *IEEE Transactions on Dependable and Secure Computing*, vol. 1, no. 1, pp. 11–33, Mar. 2004.
- [4] X. Jin, "Fault tolerant finite-time leader-follower formation control for autonomous surface vessels with los range and angle constraints," *Automatica*, vol. 68, pp. 228–236, Jun. 2016.
- [5] R. J. Patton, "Fault-tolerant control," in *Encyclopedia of Systems and Control*, J. Baillieul and T. Samad, Eds., Jul. 2015, pp. 201–213.

- [6] H. Yang, B. Jiang, H. H. T. Liu, H. Yang, and Q. Zhang, "Attitude synchronization for multiple 3-DoF helicopters with actuator faults," *IEEE/ASME Transactions on Mechatronics*, vol. 24, no. 2, pp. 597–608, Apr. 2019.
- [7] K. Zhou and Z. Ren, "A new controller architecture for high performance, robust, and fault-tolerant control," *IEEE Transactions on Automatic Control*, vol. 46, no. 1, pp. 1613–1618, Oct. 2001.
- [8] M. Blanke, M. Kinnaert, M. Staroswiecki, and J. Lunze, *Diagnosis and Fault-Tolerant Control*, 2nd ed. Heidelberg, Germany: Springer-Verlag Berlin Heidelberg, 2006.
- [9] S. Yin, H. Gao, J. Qiu, and O. Kaynak, "Descriptor reduced-order sliding mode observers design for switched systems with sensor and actuator faults," *Automatica*, vol. 76, pp. 282–292, Feb. 2017.
- [10] H. Hamadi, B. Lussier, I. Fantoni, C. Francis, and H. Shraim, "Comparative study of self tuning, adaptive and multiplexing fitc strategies for successive failures in an octorotor uav," *Robotics and Autonomous Systems*, vol. 133, Nov. 2020.
- [11] J. Jiang and Y. Zhang, "Accepting performance degradation in fault-tolerant control system design," *IEEE Transactions on Control Systems Technology*, vol. 14, no. 2, pp. 156–172, Mar. 2006.
- [12] Z. Wang, L. Liu, Y. Wu, and H. Zhang, "Optimal fault-tolerant control for discrete-time nonlinear strict-feedback systems based on adaptive critic design," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 29, no. 6, pp. 2179–2191, Jun. 2018.
- [13] I. Ahmed, H. Khorasgani, and G. Biswas, "Comparison of model predictive and reinforcement learning methods for fault tolerant control," *IFAC-Papers OnLine*, vol. 51, no. 24, pp. 233–240, Dec. 2018.
- [14] H. Zhang, K. Zhang, Y. Cai, and J. Han, "Adaptive fuzzy fault-tolerant tracking control for partially unknown systems with actuator faults via integral reinforcement learning method," *IEEE Transactions on Fuzzy Systems*, vol. 10, no. 7, pp. 1986–1998, 2019.
- [15] W. Shi, S. Song, C. Wu, and C. L. P. Chen, "Multi pseudo q-learning-based deterministic policy gradient for tracking control of autonomous underwater vehicles," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 30, no. 12, pp. 3534–3546, Dec. 2019.
- [16] D. Mguni, "Cutting your losses: Learning fault-tolerant control and optimal stopping under adverse risk," *arXiv:1902.05045*, 2019.
- [17] Q. Zhang, W. Pan, and V. Reppa, "Model-reference reinforcement learning control of autonomous surface vehicles," in *Proc. of 59th IEEE Conference on Decision and Control*, Jeju, Korea (South), Dec. 2020.
- [18] —, "Model-reference reinforcement learning for collision-free tracking control of autonomous surface vehicles," *IEEE Transactions on Intelligent Transportation Systems*, 2021, (Early access, doi:10.1109/TITS.2021.3086033).
- [19] T. I. Fossen, *Handbook of Marine Craft Hydrodynamics and Motion Control*. John Wiley & Sons, Inc., 2011.
- [20] R. Skjetne, T. I. Fossen, and P. V. Kokotović, "Adaptive maneuvering, with experiments, for a model ship in a marine control laboratory," *Mathematics of Operations Research*, vol. 41, pp. 289–298, 2005.
- [21] V. Reppa, M. M. Polycarpou, and C. G. Panayiotou, "Decentralized isolation of multiple sensor faults in large-scale interconnected nonlinear systems," *IEEE Transactions on Automatic Control*, vol. 60, no. 6, pp. 1582–1596, Mar. 2015.
- [22] —, *Sensor Fault Diagnosis*. Now Foundations and Trends, 2016.
- [23] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Proc. of the 35th International Conference on Machine Learning*, vol. 80, Stockholm, Sweden, Jul. 2018, pp. 1861–1870.
- [24] G. E. Dahl, T. N. Sainath, and G. E. Hinton, "Improving deep neural networks for lvcpr using rectified linear units and dropout," in *Proc. of 2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, Vancouver, BC, Canada, May 2013.
- [25] V. Mnih, K. Kavukcuoglu, D. Silver, and et. al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, pp. 529–533, Feb. 2015.
- [26] S. Fujimoto, H. van Hoof, and D. Meger, "Addressing function approximation error in actor-critic methods," in *Proc. of the 35th International Conference on Machine Learning Conference*, vol. 80, Stockholm, Sweden, Jul. 2018, pp. 1587–1596.