

Control of a Realistic Wave Energy Converter Model using Least-Squares Policy Iteration

Enrico Anderlini, David I. M. Forehand, Elva Bannon, and Mohammad Abusara

Abstract—An algorithm has been developed for the resistive control of a non-linear model of a wave energy converter using least-squares policy iteration, which incorporates function approximation, with tabular and radial basis functions being used as features. With this method, the controller learns the optimal PTO damping coefficient in each sea state for the maximization of the mean generated power. The performance of the algorithm is assessed against two on-line reinforcement learning schemes: Q-learning and SARSA. In both regular and irregular waves, least-squares policy iteration outperforms the other strategies, especially when starting from unfavourable conditions for learning. Similar performance is observed for both basis functions, with a smaller number of radial basis functions underfitting the Q-function. The shorter learning time is fundamental for a practical application on a real wave energy converter. Furthermore, this work shows that least-squares policy iteration is able to maximize the energy absorption of a wave energy converter despite strongly non-linear effects due to its model-free nature, which removes the influence of modelling errors. Additionally, the floater geometry has been changed during a simulation to show that reinforcement learning control is able to adapt to variations in the system dynamics.

Index Terms—Wave energy converter (WEC), resistive control, reinforcement learning (RL), function approximation, radial basis function (RBF).

I. INTRODUCTION

WAVE energy has the potential to become a significant contributor to the future energy mix thanks to a resource of up to 2.1 TW of power worldwide [1], with a consequent reduction in greenhouse gas emissions. Nevertheless, wave energy converters (WECs) are not economically viable yet, despite numerous designs having been proposed over the years. A review of some of the most promising, recent technologies can be found in [2]. The design of an effective control scheme can considerably reduce the levelised cost of energy associated with WECs, since it can bring about a gain in energy absorption with little additional hardware costs.

Since the 1970s, multiple control strategies have been studied so as to maximise energy extraction. Thorough reviews of the topic can be found in [3] for the initial analyses and

in [4] for more recent developments. By achieving resonance between the device and the incident waves, complex-conjugate control would result theoretically in optimal power absorption [3]. However, this is infeasible in practice because of the resulting large motions of the WEC in energetic sea states and the associated high loads. Hence, alternative control algorithms have been developed, which limit the motions, forces and power ratings of the device [4]. It is possible to differentiate between two main types of control schemes: time-averaged and real-time.

Time-averaged strategies assume stationary wave conditions over a prescribed time, over which a constant, optimal control setting is used [5]. In the case of reactive control, this is represented by the combination of damping and stiffness coefficients of the Power Take-Off (PTO) unit that maximise energy generation in each sea state. A specific case occurs for zero stiffness: resistive control. The optimal values are found through preliminary simulations, which can constrain the force and displacement, and then stored in a look-up table. Whereas time-averaged schemes may be less efficient than real-time strategies, their computational cost is lower.

Real-time algorithms consist of applying an optimal control action at every time instant such that it is expected to maximize energy generation over a short (in the order of one wave cycle) future time horizon [4]. Examples are latching [6], declutching [7], simple-but-effective [8], and model predictive control [9]–[11]. Whereas simple-but-effective control is computationally light, since it relies on classical closed-loop controllers, non-linear model predictive control and latching and declutching control based on Pontryagin's principle can present a high computational cost associated with their real-time optimization. Solutions to reduce the computational costs are to use a linear model in model predictive control and a moving window in latching and declutching control, as for instance proposed by [12]. Furthermore, the control operation is strongly affected by the accuracy of the wave excitation force forecast.

As the performance of all aforementioned control schemes depends on the quality of the model of the device dynamics, modelling errors can decrease the generated power as well as cause damage to the machines if the physical limits are exceeded in practice. Whereas a hierarchical robust controller has been used to decrease the sensitivity of simple-but-effective control to modelling errors and non-linear effects [13], the other control strategies are negatively affected by modelling errors. An alternative approach to robust control based on fuzzy logic has been proposed by [14]. Similarly, the authors proposed, in a previous study, the application of an alternative strategy, reinforcement learning (RL), to the control of WECs. This scheme does not rely on a model of the WEC dynamics to

Manuscript received November, 2016.

This work was supported partly by the Energy Technologies Institute and the Research Councils Energy Programme (grant EP/J500847/), partly by the Engineering and Physical Sciences Research Council (grant EP/J500847/1), and partly by Wave Energy Scotland.

E. Anderlini is with the Industrial Doctoral Centre in Offshore Renewable Energy (IDCORE), Edinburgh, EH9 3JL, UK (e-mail: E.Anderlini@ed.ac.uk).

D. Forehand is with the Institute of Energy Systems, University of Edinburgh, Edinburgh, EH9 3DW, UK (e-mail: D.Forehand@ed.ac.uk).

E. Bannon is with Wave Energy Scotland, Highlands and Islands Enterprise, Inverness, IV2 5NA, UK (e-mail: elva.bannon@hient.co.uk).

M. Abusara is with the College of Engineering, University of Exeter, Penryn, Cornwall, TR10 9FE, UK (e-mail: M.Abusara@exeter.ac.uk).

obtain the control action and is thus able to adapt to changes in the system response due to ageing, e.g. the build-up of marine biofouling. Although the algorithm described in [15] can be implemented on an actual device, a test-case numerical study was run in that paper using a linear model of a point absorber, a well known WEC technology which comprises of a float whose size is small compared with the characteristic wavelength [2].

In this article, RL is applied to the control of a realistic, non-linear model of a WEC, whose accuracy was validated by the developers against measurements on a prototype device [16]. In the absence of costly experimental measurements, this study enables the assessment of the convergence behaviour of RL when non-linear effects are important. As in [16], only damping, or resistive, control is analysed. This is known to be inferior in performance to fully reactive control, which is in fact treated by most other real-time control schemes, as described in [4]. However, the practical implementation of resistive control is simple. Additionally, the optimization of only the damping coefficient results in a simpler framework to demonstrate the applicability of RL for the control for WECs, which can then be extended to the treatment of combined damping and stiffness control. In particular, here resistive control is implemented using least-squares policy iteration (LSPI), an efficient RL algorithm [17]. Furthermore, its performance is assessed against two simpler RL algorithms, SARSA and Q-learning [18]. In addition, the effectiveness of function approximation in reducing the learning time has been assessed using LSPI with radial basis functions (RBFs) [17].

II. RESISTIVE CONTROL OF THE SEABASED WEC

A. System Description

The Seabased device is a point absorber with a direct-drive PTO system. The development and testing of a number of full-scale prototypes at Uppsala University is well described in the literature [19]–[23]. The version studied in [16] is analysed in this article, although the generator is now connected to the electrical grid.

Figure 1 shows a diagram of the device, which is inspired by [16]. A small float, excited by incident waves, drives a linear, permanent-magnet generator along vertical rails. The two bodies are connected by a mooring line. When the distance between the float and the translator decreases, the mooring line goes slack and the translator is pulled downwards by a dedicated spring. Additionally, springs at the upper and lower end stops prevent the translator from breaking the casing in large waves. The motion of the magnet induces electrical current in the coils wound around the stator. Power absorption is controlled through a power electronic converter by setting the stator current I_s to be proportional to the velocity of the translator. A second power electronic converter controls the voltage across the capacitor between the converters by setting the grid current. The wave elevation ζ is measured through a wave buoy sited 80 m from the prototype at the Lysekil wave energy research site [21].

In Figure 1, the same naming convention as in [16] is held, with the values of the variables quantities being given in

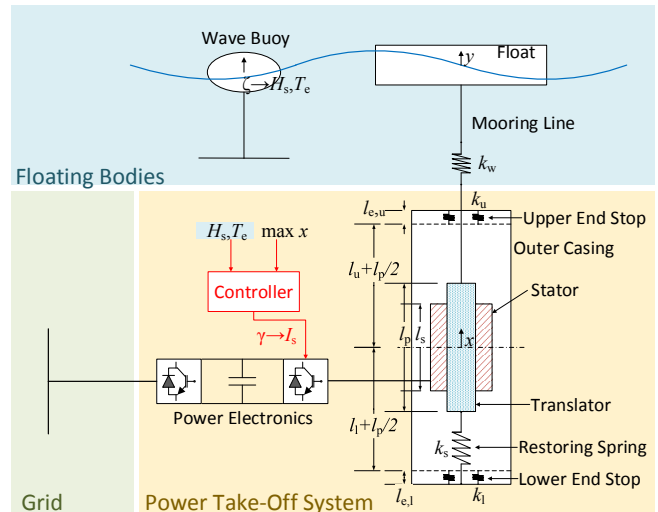


Fig. 1. Diagram of the prototype Seabased WEC.

Table I in [16]. In addition, $l_{e,u} = 0.25$ m and $l_{e,l} = 0.14$ m are a measure of the end stops length, as given in [21].

B. Mathematical Model

A weakly non-linear mathematical model of the system dynamics has been developed by [16]. Although the float is free to move in all directions in reality, only the heave degree of freedom is analysed because the influence of the other motions is negligible [24]. Defining y and x as the float and translator displacement respectively, the motions of the two bodies are described by the following system of equations

$$(m_b + m_\infty) \ddot{y}(t) = F_e(t) - F_r(t) - F_h(t) - F_w(t), \quad (1)$$

$$m_p \ddot{x}(t) = F_w(t) - F_{em}(t) - F_s(t) + F_u(t) + F_l(t), \quad (2)$$

where t indicates time, m_b and m_p the mass of the float and piston respectively, and m_∞ the added mass of the float. F_e is the incident and diffracted wave excitation force, F_r part of the float radiation force, F_h the hydrostatic restoring force, F_w the tension in the wire connecting the float to the translator, F_{em} the electromotive force, F_s the force of the restoring spring in Fig. 1, F_u and F_l the spring force of the upper and lower end stops, respectively. The non-linearities are associated with F_w , where compression effects are ignored, F_u and F_l , which are activated only if the end stops are reached, and F_{em} , which depends on the exposure of the translator to the stator.

The electromotive, or control, force F_{em} is discussed in the next section. The hydrostatic force is calculated as

$$F_h(t) = \rho g S_w y(t), \quad (3)$$

where $\rho = 1025$ kg/m³ is the seawater density, $g = 9.81$ m/s² the gravitational acceleration and S_w the float waterplane area.

The F_r part of radiation force is approximated through a state-space system so as to reduce its computational cost [25]

$$\dot{\mathbf{x}}_{ss}(t) = \mathbf{A}_{ss} \mathbf{x}_{ss}(t) + \mathbf{B}_{ss} \dot{y}(t) \quad (4)$$

$$F_r(t) \approx \mathbf{C}_{ss} \mathbf{x}_{ss}(t). \quad (5)$$

The matrices A_{ss} , B_{ss} and C_{ss} are calculated with frequency-domain system identification, as described in [25]. F_r is given by the convolution of the product of the radiation impulse response function and the float velocity [26]. Furthermore, the radiation impedance function has been computed using the commercial software WAMIT for the vertical cylinder geometry described in [16] for circular wave frequency values ranging from 0 rad/s to 10 rad/s in steps of 0.005 rad/s. Furthermore, m_∞ has been calculated for the infinite wave frequency case.

Similarly, the wave excitation force F_e is given by the convolution integral of the product of the wave elevation at the float location and the diffraction impulse response function [26], which has also been computed using WAMIT. In irregular waves, a number of wave components have been superimposed to obtain ζ . The amplitude of each wave has been obtained from a wave spectrum, sampled at a circular wave frequency step of 0.005 rad/s in order to prevent repeating the wave trace within a 15-minute window, as the value is less than the Nyquist frequency [27]. In order to obtain longer time series, individual wave traces generated using a different seed to the random number generator have been joined. The connections are smoothed using a 200-point filter over the last and first 20 s of each wave trace.

The force of the spring connected to the translator is expressed as [16]

$$F_s = F_0 + k_s x, \quad (6)$$

where F_0 is a static force due to precharging, and k_s is the spring stiffness. Ignoring compression effects, the wire force is given by [16]

$$F_w = \begin{cases} -k_w(y - x) & \text{if } y > x \\ 0 & \text{else} \end{cases}, \quad (7)$$

with k_w being the wire stiffness. Similarly, the forces due to the upper and lower end stops are given by [16]

$$F_u = \begin{cases} -k_u(x - l_u) & \text{if } x > l_u \\ 0 & \text{else} \end{cases}, \quad (8)$$

$$F_l = \begin{cases} -k_l(x + l_l) & \text{if } x < -l_l \\ 0 & \text{else} \end{cases}, \quad (9)$$

where k_u and k_l are the equivalent stiffness values of the springs in the upper and lower end stops respectively. l_u and l_l are the distance of the two end stops from the vertical midpoint of the translator at equilibrium, as shown in Fig. 1. In [16], it is possible to find the values of l_l , l_u , k_l , k_u , k_w , k_s , m_b , m_p , F_0 and S_w .

Using (3-9), Equations (1-2) have been expressed in the following non-linear state-space form

$$\dot{z}(t) = \mathbf{A}z(t) + \mathbf{B}u(t, x) + \mathbf{B}w(t) + \mathbf{B}l(t, x, y), \quad (10)$$

The state, input, noise and non-linear vectors are given by

$$z = [y \quad \dot{y} \quad x \quad \dot{x} \quad \mathbf{x}_{ss}^T]^T, \quad (11)$$

$$u = [0 \quad -F_{em}(x)]^T, \quad (12)$$

$$w = [F_e(t) \quad 0]^T, \quad (13)$$

$$l = [-F_w(x, y) \quad F_w(x, y) - F_0 + F_u(x) + F_l(x)]^T. \quad (14)$$

The state and input matrices are

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & 0 & \mathbf{0}^T \\ -\frac{\rho g S_w}{m_b + m_\infty} & 0 & 0 & 0 & -\frac{C_{ss}}{m_b + m_\infty} \\ 0 & 0 & 0 & 1 & \mathbf{0}^T \\ 0 & 0 & -\frac{k_s}{m_p} & 0 & \mathbf{0}^T \\ \mathbf{0} & \mathbf{B}_{ss} & \mathbf{0} & \mathbf{0} & \mathbf{A}_{ss} \end{bmatrix}, \quad (15)$$

$$\mathbf{B} = \begin{bmatrix} 0 & 0 \\ \frac{1}{m_b + m_\infty} & 0 \\ 0 & 0 \\ 0 & \frac{1}{m_p} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (16)$$

C. Resistive Control of the Translator

The motions of the float and the translator, and thus ultimately the power absorption of the device, can be controlled through the electrical behaviour of the generator. In particular, the electromotive force is proportional to stator current I_s and active area A_{fac}

$$F_{em} = k_\tau A_{fac}(x) I_s, \quad (17)$$

where k_τ is the generator torque constant. If the current is controlled (by power electronics) so that it is proportional to speed such as $I_s = b\dot{x}$, with b being a constant, then Equation (17) becomes

$$F_{em} = k_\tau b A_{fac}(x) \dot{x}, \text{ or } F_{em} = \gamma A_{fac}(x) \dot{x}, \quad (18)$$

where $\gamma = k_\tau b$ is the PTO damping coefficient. The active area, i.e. the overlap between stator and translator, is given by

$$A_{fac} = \begin{cases} 0 & \text{if } |x| \geq 0.5(l_p + l_s) \\ 1 & \text{if } |x| \leq 0.5(l_p - l_s) \\ [0.5(l_p + l_s) - |x|] / l_s & \text{else} \end{cases}$$

with l_p and l_s being given in [16]. The generated power is computed as

$$P(t) = F_{em}(t) \dot{x}(t). \quad (19)$$

RL is employed to find the optimal PTO damping coefficient, γ , in each sea state for the maximization of the energy generation. The values of γ are assumed to be limited to 0-100 kNs/m. The upper limit corresponds approximately to the case of no load resistance in the experimental setting in [16].

D. Simulation System

Equation 10 has been discretized with a fourth-order Runge-Kutta scheme [28], and solved with a time step of 0.002 s. The controller is implemented as in (18). The workflow of the program is similar to that described in [15], and it is summarised in Fig. 2 with a block diagram.

III. REINFORCEMENT LEARNING

In the RL framework [29], an agent takes an action a in state s , landing in a new state s' while observing a reward r . A Markov decision process is used to model the action selection depending on the value function $Q(s, a)$, which represents an estimate of the future reward. With time, the agent learns an optimal policy, π , that maximizes the total reward.

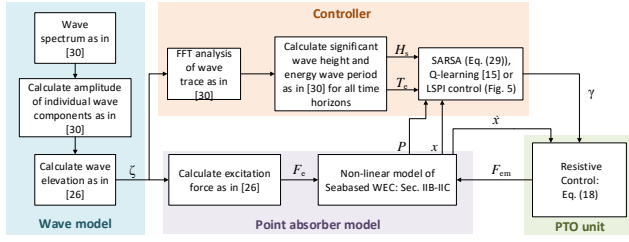


Fig. 2. Workflow diagram of the computer program used to simulate the Seabased WEC.

A. Least-Squares Policy Iteration

LSPI is a powerful, off-line, off-policy RL algorithm. Whilst still being model-free, the method, developed by [17], presents an efficient use of the samples (s, a, s', r) , which results in a smaller learning time as compared with other strategies, such as Q-learning and SARSA. Additionally, it automatically incorporates function approximation for the action-value function. This means that the scheme is able to generalize for unseen states, thus further shrinking the convergence time. In particular, a linear architecture is used for the approximation of Q due to its simple implementation and ease of debugging and feature engineering [17]. In matrix notation, this is expressed as [18]

$$Q(s, a) \approx \phi(s)^T \Theta_{:,a}, \quad (20)$$

where Θ is the weight matrix and ϕ is the vector of arbitrary, linearly independent, usually non-linear basis functions, or features. $\Theta_{:,a}$ indicates the a^{th} column of Θ , with Θ having $|\mathcal{A}|$ columns, where \mathcal{A} is the action space. Θ and ϕ have $J \ll |\mathcal{S}|$ rows, with \mathcal{S} indicating the state space. Here, two basis functions types are used: tabular and radial. The tabular representation is the simplest and consists in assigning a separate weight for each state-action pair [18]. Hence, for discrete states, this corresponds to the exact representation $Q(s, a)$, although its size is equal to the state-action space ($J = |\mathcal{S}|$). Conversely, in RBFs, the feature activation decays continuously away from the state-action pair where the RBF is centred, s_j for RBF j , spanning many discrete states [18]

$$\phi_j(s) = \exp\left(-\frac{\|s - s_j\|^2}{2\mu_j}\right), \quad (21)$$

where μ_j indicates the bandwidth of RBF j . RBFs are shown graphically in Fig. 3.

LSPI consists of two main stages: policy evaluation (the critic) and policy improvement (the actor) [17]. LSPI is defined as off-line because the algorithm is trained using samples that have been previously recorded from observations of the environment. The algorithm is summarized in Fig. 4. The discount factor is set here to $\gamma_d = 0.95$. The values of the weight matrix in (20) can be computed from

$$\tilde{\mathbf{A}}\Theta_{:,a} = \tilde{\mathbf{b}} \quad (22)$$

for each action, where the tilde indicates a learned variable. The reader is referred to [17] for a full derivation of the equations for $\tilde{\mathbf{A}}$ and $\tilde{\mathbf{b}}$, which are obtained from the least-squares fixed-point approximation.

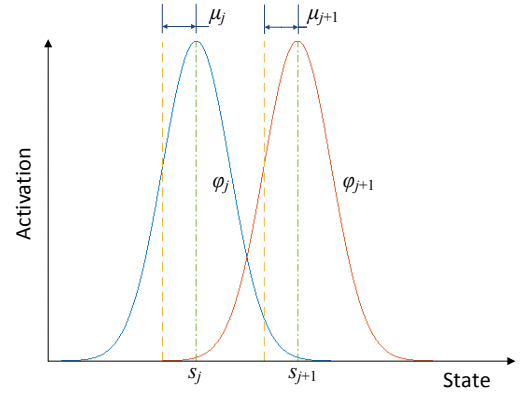


Fig. 3. Activation function of the RBFs as per Eq. (21).

input: W : set of samples (s, a, r, s')
 γ_d : discount factor
 $\delta = 10^{-3}$: stopping criterion
 π_0 : initial policy, given as $\theta_0 = \mathbf{0}$
 π : policy, or exploration strategy

- $\theta' \leftarrow \theta_0$
- while $\|\theta - \theta'\| \geq \delta$:
 - $\theta \leftarrow \theta'$
 - $\tilde{\mathbf{A}} \leftarrow \mathbf{0}$ ($J \times J$ matrix)
 - $\tilde{\mathbf{b}} \leftarrow \mathbf{0}$ (J vector)
 - for each $(s, a) \in W$:
 - $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}} + \phi(s)(\phi(s) - \gamma_a \phi(s', \pi(s')))^T$
 - $\tilde{\mathbf{b}} \leftarrow \tilde{\mathbf{b}} + \phi(s)r$
 - $\theta'_{:,a} \leftarrow \tilde{\mathbf{A}}^{-1}\tilde{\mathbf{b}}$
- return θ

Fig. 4. LSPI algorithm, adapted from [17].

B. Application of LSPI to the control of WECs

Employing a time-averaged approach, at the start of each time-averaging period, or time horizon with duration H_{RL} , an action, which consists in a step change in the PTO damping coefficient, is selected following the current policy. The state is a combination of γ and the sea state, as given by the significant wave height H_s and energy wave period T_e [30]. Holding γ constant during H_{RL} , the reward is obtained as a function of the mean generated power, P_{avg} . The selection of a new action results in a new state, and the sample (s, a, s', r) is added to the sample set W . After the collection of N_s samples, the policy is updated using the LSPI algorithm in Fig. 4. In the following sections it is possible to find an accurate description of the state and action spaces, the reward function and the exploration strategy.

1) *State Space*: Similarly to [15], the discrete state-space can be expressed as

$$\mathcal{S} = \left\{ s | s_{i,l,m} = (H_{s,i}, T_{e,l}, \gamma_m), \begin{matrix} i = 1 : I, \\ l = 1 : L, \\ m = 1 : M \end{matrix} \right\}. \quad (23)$$

As described in Sec. III-A, LSPI incorporates linear function approximation. With the tabular approach, $J = ILM$, i.e. there is an entry in Θ for each state-action value, or the Q -table is exact. With RBFs, a smaller number of values can be used. In fact, for the control of WECs, a hybrid approach is

TABLE I
DISTANCE BETWEEN KERNELS, BANDWIDTH AND NUMBER OF KERNELS USED IN THE STUDY OF LSPI WITH RBFs.

δ_c (kNs/m)	μ (kNs/m)	M
10	10	10
10	20	10
20	10	5
20	20	5
20	40	5

used, where discrete sea states are still employed, while RBFs approximate the control variable. I and L are determined from the wave data at the deployment site, with steps of 1 m and 1 s being common for H_s and T_e , respectively [30].

For the tabular approach, $M = 11$ has been selected, with γ ranging from 0 to 100 kNs/m in steps of 10 kNs/m. For function approximation, 5 cases have been considered in order to study the influence of the number of kernels, or centres, and bandwidth on the learning behaviour of LSPI with RBFs. In Table I, it is possible to see the distance between kernels $\delta_c = s_j - s_{j-1}$ and bandwidth μ for each case as given in (21). The first kernel is always sited at $\gamma = 0$ kNs/m.

2) *Action Space*: The action space is defined as

$$\mathcal{A} = \{a | (-\Delta\gamma, 0, +\Delta\gamma)\}, \quad (24)$$

where $\Delta\gamma = \gamma_{m+1} - \gamma_m$. However, the states corresponding to the maximum and minimum of the PTO damping coefficient are constrained to two actions in order to avoid going beyond the RL state space limits.

3) *Reward Function*: The same reward function as in [15] is employed and the reader is referred there for a more detailed explanation. A penalty $p = -1$ is returned if the constraints $\max(x) > l_u + l_{e,u}$ or $\min(x) < -(l_l + l_{l,u})$ are exceeded during the time horizon h . This enables the algorithm to learn to avoid actions that will result in possible damage to the device. Hence, the reward function is expressed as

$$r = \begin{cases} \left[\frac{\mathbf{m}(s_h)}{\max_{s''=o:p} \mathbf{m}(s'')} \right]^u & \text{if constraints met} \\ p & \text{otherwise} \end{cases} \quad (25)$$

The entries of the vector \mathbf{m} , whose size is equal to the total number of discrete states $|\mathcal{S}|$ (valid for both tabular and RBF methods), correspond to the average of up to 10 values of P_{avg}/H_s^2 that are stored for each discrete state, with older values being overwritten by new ones once 10 values are registered. The indices o and p ensure that the maximization in (25) is performed only over the values of \mathbf{m} corresponding to the current sea state, as given by H_s and T_e .

The power u must be an odd number to prevent rectifying negative power values (in the case of reactive control). The higher the value, the closer the cost function is to returning 1 for the optimal control variable and 0 for all other settings in each sea state. Here, $u = 25$ has been used in order to aid convergence in irregular waves.

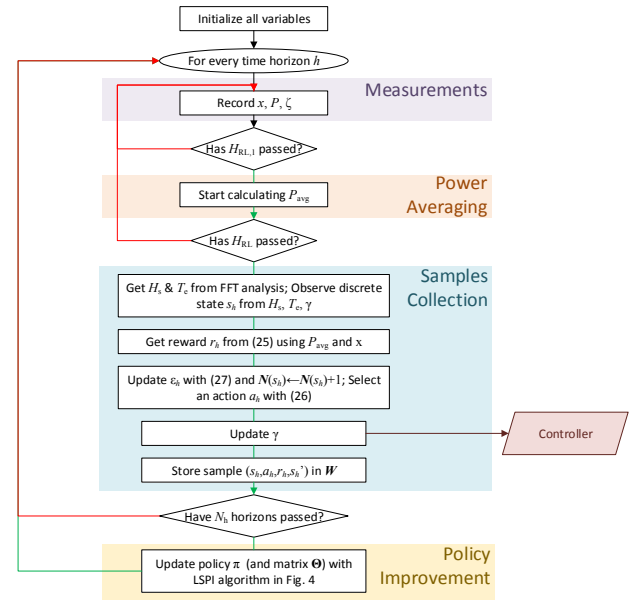


Fig. 5. Flowchart of the LSPI algorithm for the resistive control of the stator.

4) *Exploration Strategy*: An ϵ -greedy policy [29] selects the action at the start of each time horizon h

$$a = \begin{cases} \arg \max_{a' \in \mathcal{A}} Q(s_h, a') & \text{with probability } 1 - \epsilon_h \\ \text{random action} & \text{with probability } \epsilon_h \end{cases}, \quad (26)$$

where ϵ_h is the exploration rate. $\max_{a' \in \mathcal{A}} Q(s_h, a')$ represents the maximum action-value (i.e. a measure of the expected reward) for the current state over all actions, with the action-value function being given by the mapping in (20). This term represents the selection of the action that results in maximum expected total reward starting from the current state.

Greater exploration is desired at the start of RL control, while the greedy action, i.e. such that it maximises the value function, is preferred as the learned policy improves. Thus, the exploration rate is obtained as

$$\epsilon_h = \begin{cases} \epsilon_0 & \text{if } N(s_h) \leq N_\epsilon \\ \epsilon_0 / \sqrt{N(s_h) - N_\epsilon} & \text{if } N(s_h) > N_\epsilon \end{cases}, \quad (27)$$

with $N(s_h)$ indicating the number of visits to the current discrete state (hence, valid for both tabular and RBF approaches). $N_\epsilon = 5$ is the minimum number of encounters for random exploration, and the initial exploration rate is set to $\epsilon_0 = 0.5$.

C. Algorithm

The proposed LSPI algorithm for the resistive control of the stator can be seen in Fig. 5. After initializing all variables, the algorithm is run continuously until the device is disconnected, e.g. due to maintenance. During each time horizon h , the policy is applied in order to select a suitable action based on the encountered sea state, mean generated power and maximum translator displacement. Furthermore, at the end of each horizon, the current state, action, next state and reward are sampled and added to \mathcal{W} . Due to the finite memory of

the controller computer, a specified number of samples can be stored, say 10^6 . Therefore, new samples will be stored only if they have not been recorded before, with a difference greater than 10^{-3} being acceptable for the reward. Once the memory limit is reached, older values will need to be overwritten, ensuring the sample range is broad, i.e. accounting for the different sea states and values of the PTO damping coefficient.

As shown in Fig. 5, the policy is improved using the LSPI algorithm in Fig. 4 every $N_h = 40$ time horizons. This operation can be performed off-line on separate computing cores so as to reduce the computational effort and ensure the on-line implementation is feasible.

A time horizon duration $H_{RL} = 10T$ has been chosen in regular waves, with T being the wave period, while $H_{RL} = 150$ s in the analysed irregular waves because a JONSWAP spectrum is used. This selection is based on a compromise between a fast response and a stable algorithm. Irregular waves in particular require a longer duration of the power averaging process due to their stochastic nature. If a wider-banded wave spectrum is adopted, the horizon length should be increased. Additionally, this process is started only after $H_{RL,1} = 0.4H_{RL}$ so as to remove the influence of the transient effects associated with the change in load resistance.

D. Q-learning and SARSA

The performance of LSPI in the control of the WEC is compared with Q-learning and SARSA. SARSA, which stands for *state-action-reward-state-action*, and Q-learning are on-line schemes that rely on discrete states and actions [29]. Hence, at each step, they update the Q-table with the following equations [29]

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma_d \max_{b \in \mathcal{A}} Q(s', b) - Q(s, a) \right], \quad (28)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \left[r + \gamma_d Q(s', a') - Q(s, a) \right], \quad (29)$$

for Q-learning and SARSA respectively. α is the learning rate and a' the action applied in the future state. Hence, it is clear that the main difference between the two algorithms is that while SARSA is on-policy, i.e. it updates the value function based on the policy it will follow, Q-learning is off-policy, i.e. the update is based on the maximum possible Q-value in the new state [29]. The application of Q-learning to the resistive control of WECs is described in a previous publication [15], with SARSA presenting an almost identical implementation. Hence, the reader is referred to [15] for details. The initial learning rate is set here to 0.4.

IV. SIMULATION RESULTS

A. Regular Waves

The behaviour of SARSA, Q-learning and LSPI has been assessed against the optimal PTO damping coefficient, which has been calculated using the Matlab optimization function *fmincon* in each sea state. This is to provide a benchmark of the control variable that results in the maximum mean generated power.

Regular waves of unit amplitude and a wave period of 6 s have been analysed first, with a wave trace lasting 3

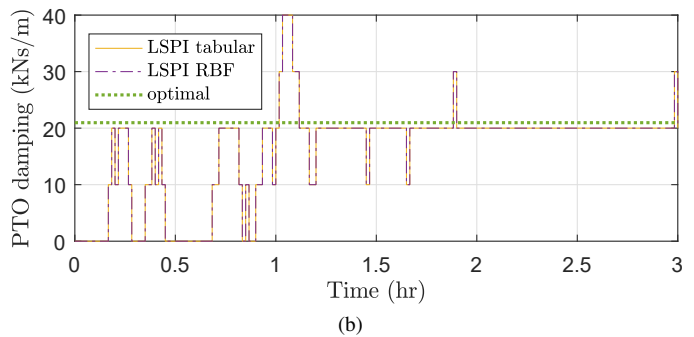
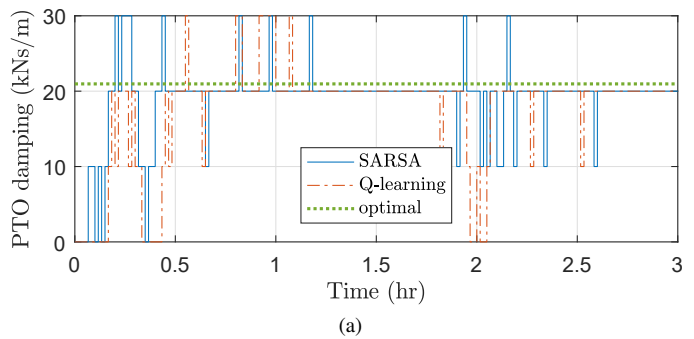


Fig. 6. PTO damping coefficient selected by different RL control strategies as compared with the optimal value in regular waves with unit amplitude and $T = 6$ s starting from $\gamma = 0$ kNs/m.

hours. Two different starting points have been selected, namely $\gamma = 0$ and $\gamma = 100$ kNs/m, as shown in Fig. 6 and Fig. 7, respectively. For the RBFs, $\delta_c = 10$ kNs/m and $\mu = 10$ kNs/m, i.e. an almost tabular approach has been used. For each figure, the same seed number has been set to the random number generator for all algorithms, selecting a particularly unfavourable number for Fig. 7 in order to assess the convergence properties under difficult conditions.

In Fig. 8, it is possible to see the behaviour of the LSPI algorithm for the RBF settings in Table I, when the starting value of the PTO coefficient is $\gamma = 100$ kNs/m. For all runs, the same seed values is used as in Fig. 7. A longer wave trace lasting 4 hours is employed.

The mean generated power corresponding to the run with LSPI with RBFs and $\delta_c = 10$ kNs/m and $\mu = 10$ kNs/m in Fig. 7b and Fig. 8 is plotted in Fig. 9.

B. Irregular Waves

In irregular waves, an 8-hour long wave trace with $H_s = 2$ m and $T_e = 6$ s with a JONSWAP spectrum [30] has been analysed, typical of the Lysekil testing site [31]. In Fig. 10a and Fig. 10b, the learning behaviours of the three control algorithms are shown, with the same setting being used for LSPI with RBFs as in Fig. 7 throughout this section. The difference in mean generated power between LSPI with RBFs and the optimal control setting is shown in Fig. 10c.

Nevertheless, real sea states actually last between 0.5 to 6 hours [30]. Therefore, in order to prove that RL is able to deal with changing sea states, the control is tested in an additional 12-hour-long wave trace composed of the alternation of two sea states, so that $I = L = 2$. Both have a JONSWAP

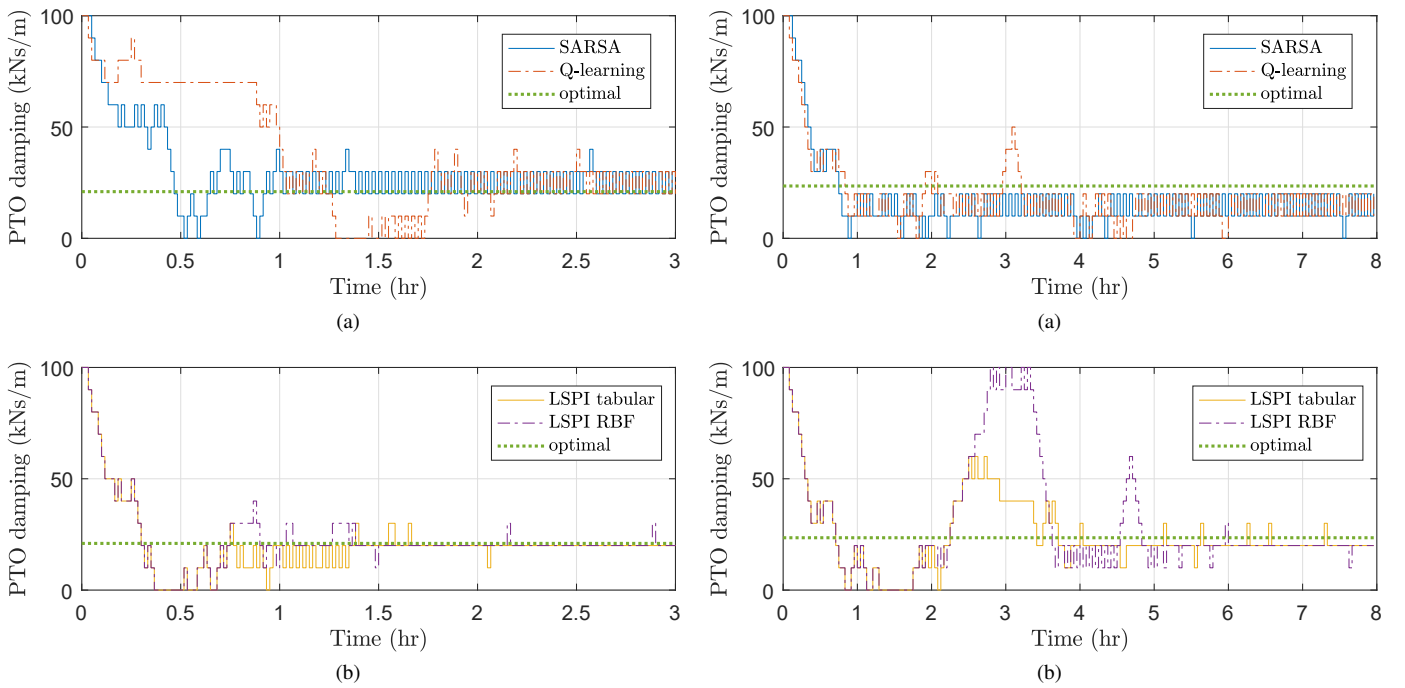


Fig. 7. PTO damping coefficient selected by different RL control strategies as compared with the optimal value in regular waves with unit amplitude and $T = 6$ s starting from $\gamma = 100$ kNs/m.

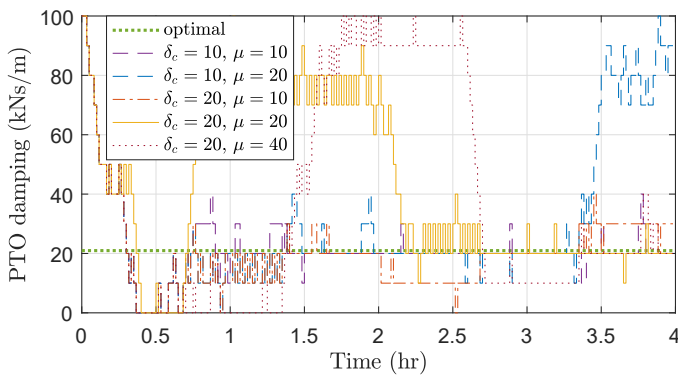


Fig. 8. PTO damping coefficient selected by the LSPI algorithm with different RBF settings in regular waves with unit amplitude and $T = 6$ s. The values of δ_c and μ are in kNs/m.

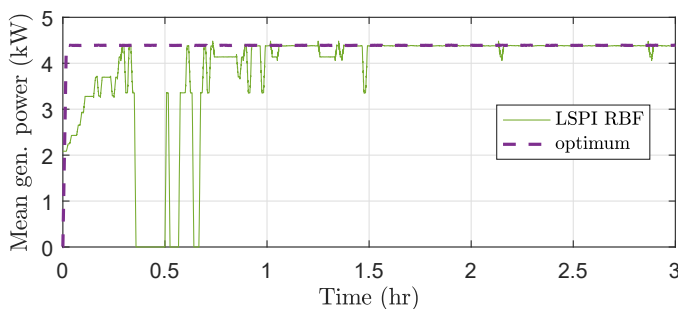


Fig. 9. Mean generated power for the run with LSPI with RBFs and $\delta_c = 10$ kNs/m and $\mu = 10$ kNs/m in Fig. 7b and Fig. 8.

spectrum and last for two hours before changing. The first one corresponds to $H_s = 2$ m and $T_e = 5$ s, while the second one has $H_s = 1$ m and $T_e = 6$ s. Fig. 11a and Fig. 11b show

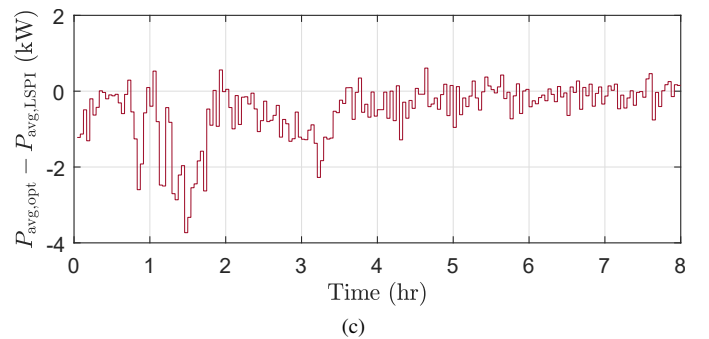


Fig. 10. PTO damping coefficient selected by different RL control strategies as compared with the optimal value in irregular waves with $H_s = 2$ m and $T_e = 6$ s and a JONSWAP spectrum starting from $\gamma = 100$ kNs/m (a-b). (c) shows the difference in the mean generated power for the optimum ($P_{avg,opt}$) and the case of LSPI with RBFs ($P_{avg,LSPI}$).

the learning behaviour of the three RL algorithms. In fig. 11c, the difference in mean power between LSPI with RBFs and the optimal control setting in each sea state can be seen.

Furthermore, although RL is expected to result in adaptive control, as it is model-independent [32], this was not proven in the previous work on the control of WECs [15]. Hence, a simple example is treated here to show the adaptivity of RL to possible marine growth effects. Bio-fouling is expected to affect the dynamics of the system mainly through an increase in its inertia and especially drag force. However, in this simple model, the viscous drag force is not considered. Hence, we treat the case of a sudden increase in the radius and draught of the floater to 1.75 m and 0.5 m, respectively (from 1.5 m and 0.4 m, respectively, in [16]). These values have been assumed, as they result in a significant change in the optimal damping coefficient in the analysed sea state. A full sensitivity analysis of the power absorption and control of the device to the variations in floater design as well as a realistic treatment of marine growth effects go beyond the scope this study. The

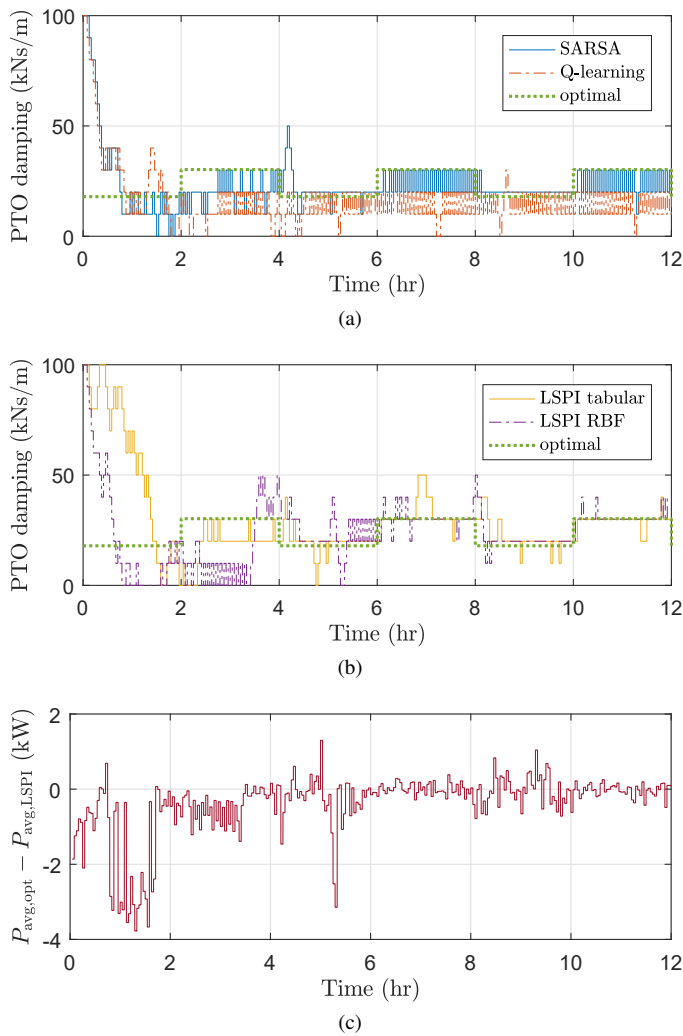


Fig. 11. PTO damping coefficient selected by different RL control strategies as compared with the optimal value in irregular waves with two alternating sea states (JONSWAP spectra with $H_s = 2$ m and $T_e = 5$ s, and $H_s = 1$ m and $T_e = 6$ s) starting from $\gamma = 100$ kNs/m (a-b). (c) shows the difference in the mean generated power for the optimum ($P_{avg,opt}$) and the case of LSPI with RBFs ($P_{avg,LSPI}$).

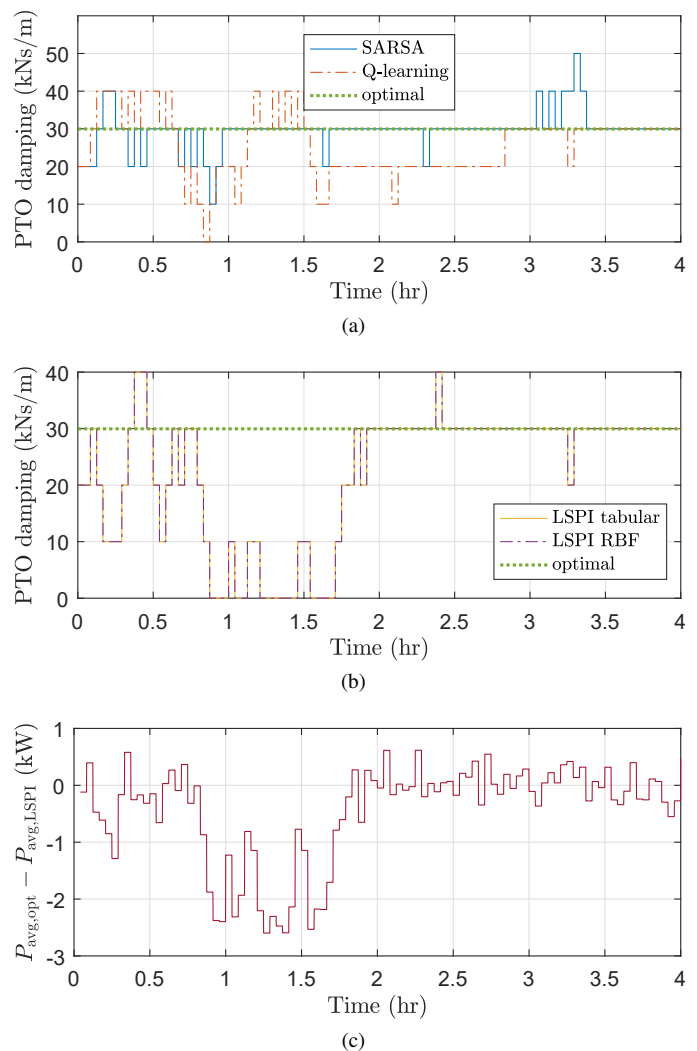


Fig. 12. PTO damping coefficient selected by different RL control strategies as compared with the optimal value for the new floater geometry in irregular waves with $H_s = 2$ m and $T_e = 6$ s and a JONSWAP spectrum. The initial conditions are set based on the final settings of Fig. 10a and Fig. 10b, respectively. (c) shows the difference in the mean generated power for the optimum ($P_{avg,opt}$) and the case of LSPI with RBFs ($P_{avg,LSPI}$).

changes in the floater design result in an increase in its mass (m_b) of 2032.7 kg, in its surface area of 2.5525 m² (note that these values are relative to the values in [16]) as well as changes in the radiation approximation state-space system, with the radiation coefficients being computed with WAMIT.

The same sea state as in Fig. 10 is used in this simple example, whereas the new geometry of the floater is employed. In particular, a simulation is initialized with the final values of Fig. 10a and Fig. 10b being set for each RL strategy. Additionally, the same values of the \mathbf{m} vector have been kept for each scheme. This corresponds to initializing the reward function to incorrect values for each discrete damping coefficient. For this reason, the exploration rate (as well as the learning rate for Q-learning and SARSA) is reinitialized with the same settings as in Sec. III-B4.

In Fig. 12a and Fig. 12b, the learning behaviours of the three control algorithms are shown. The difference in mean generated power between LSPI with RBFs and the optimal control setting is shown in Fig. 12c.

The computational time of the algorithm run at the start of each time horizon has been less than 0.06 s on an i7 processor with 16Gb RAM in all simulations run here. As this time is proportional to the number of states, if, say, 100 sea states were to be used, the computational time would increase to 0.3 s. Hence, a practical implementation is realistic, particularly considering the much longer time horizon duration.

V. DISCUSSION

A. Regular Waves

In this work, we define RL algorithms to have converged towards a policy once the same PTO damping coefficient is selected for longer than an hour. However, within the short duration of the analysed wave traces, the exploration rate does not fully decay. Hence, the definition of convergence is extended to include a maximum of up to 5 distinct deviations from the mean value of the selected γ within the one-hour period, which may be due to random actions being adopted.

In Fig. 6, it can be seen that all algorithms learn the optimal PTO damping coefficient within 2.5 hours, with subsequent wiggles, especially visible for Q-learning and SARSA, mainly due to the exploration rate not having fully decayed. This fast learning is because this is a benign case, with the optimal value of γ being very close to the starting PTO damping coefficient, thus requiring little exploration before finding the optimum. Conversely, Fig. 7 represents a more challenging scenario for the RL algorithms. In particular, SARSA and Q-learning are unable to converge to the optimal policy, and learn a suboptimal policy instead, which results in less energy absorption than the optimal policy. This problem could be solved with a slower decay in the exploration and learning rates, which would cause learning to be smoother, but also slower. This behaviour is particularly worrying in the case of extreme waves because if this oscillation occurs on the boundary of the feasible damping coefficient envelope to prevent excessive displacements, it could lead to failure. Conversely, LSPI with both tabular and radial basis functions learns the optimal policy within 2.5 hours in regular waves in Fig. 7b.

Comparing the behaviour of LSPI with tabular features and RBFs with $\delta_c = 10$ kNs/m and $\mu = 10$ kNs/m in Fig. 6b, 7b and Fig. 10b, the two approaches almost completely match, with RBFs actually resulting in a stabler behaviour in regular waves and greater exploration in irregular waves. This is expected because almost the same number of kernels as discrete states are used, with the bandwidth spanning the space between discrete states. In Fig. 4, decreasing the number of kernels was expected to result in faster learning because the RBFs are expected to generalise the shape of the Q-function for unseen states and actions [18]. In fact, this is not the case, with LSPI with RBFs with $\delta_c = 20$ kNs/m (thus half as many kernels) and $\mu = 20$ kNs/m taking longer to learn the optimal policy. Increasing the bandwidth of RBFs also augments the confusion in the controller, as the overlap between distinct RBFs is increased spanning multiple γ values, thus causing the algorithm to diverge from the optimal policy. These counter-intuitive results are believed to be due to the small number of discrete states used, with many more features being typical for standard RL problems [17]. Hence, the use of 5 or less RBFs incurs in an underfitting problem, i.e. using too coarse a model to fit the Q-function. A minimum of 10 RBFs is recommended for the control of WECs with LSPI. Additionally, setting the bandwidth to match the distance between kernels seems to provide best behaviour. Nevertheless, designing RBFs features needs care, and it is likely to be device-specific.

B. Irregular Waves

Q-learning and SARSA are similarly unable to converge towards the optimal policy in irregular waves as well, as shown in Fig. 10 and Fig. 11. Again, this is an indicator that the exploration and learning rates should be decreased more slowly for these algorithms, thus resulting in longer learning times. Conversely, LSPI with both tabular and radial basis functions is able to learn the optimal policy in less than 6 hours in each sea state, despite some wiggles owing to the exploration rate not having decayed fully yet in Fig. 11b.

In particular, the learning time is lower than the 12 hours required by Q-learning for convergence in irregular waves in [15], where a more benign linear WEC model was used for validation. This diminished convergence time is mainly due to the shorter time-averaging horizon length employed in this study and, especially, the superior capacity of LSPI to learn using a small number of observations [17]. Furthermore, as shown in Fig. 11b, LSPI is able to pick up learning in a specific sea state from where it left off the last time the controller was in that sea state. This is a fundamental consideration for a realistic application, since actual sea states usually last for a shorter time than 6 hours [30].

As the Seabased device is tested in the Skagerrak strait [16], a JONSWAP spectrum is appropriate due to its bounded, shallow-water nature [30]. However, a JONSWAP spectrum is a single-peaked spectrum with a relatively narrow frequency range [30]. This means that energy is contained mainly in a region close to the peak wave period. As a result, determining the optimal PTO damping coefficient for each sea state is simpler than for wider-banded wave spectra, such as Bretschneider or even double-peaked spectra. Although RL is expected to find the global optimum [29], the learning process would be expected to take longer if the latter spectra were used: a longer time horizon length would be necessary. In particular, a double-peaked spectrum would cause significant challenges to the convergence behaviour. This will be the focus of future studies.

Being model-free, RL is proven to be able to adapt to changes in the dynamics of the WEC in Fig. 12. Even though the reward function is initialized with the wrong values, RL is able to converge towards the optimal PTO damping coefficient with all three analysed algorithms. However, it is important to note that this is possible because the exploration rate is reset after the change of the system dynamics. Therefore, during operation of a WEC, it is necessary to reset the exploration rate after specific time intervals, say yearly, in order to pick up any possible changes in the device response.

VI. CONCLUSION

An efficient RL algorithm has been suggested for the control of a WEC, with its performance being compared with Q-learning and SARSA. In particular, a non-linear model of the dynamics of the Seabased point absorber, validated in a previous study, has been used as a test case. As expected, despite the system non-linearities, all control schemes are able to find the optimal PTO damping coefficient from a random start in regular waves because of their model-free nature. However, if the algorithms are started with particularly unfavourable conditions, only LSPI is able to converge within 2.5 hours, with higher learning and exploration rates being required for Q-learning and SARSA to converge. Unexpected results have been found in the study of RBFs as features for function approximation with LSPI: a smaller number of RBFs than discrete states does not correspond to faster learning time. This is because a very small number of discrete states has been employed, with the few RBF kernels resulting in underfitting. Hence, although RBFs should be preferred over

tabular features as they presented a stabler behaviour, their number should be high enough to prevent underfitting, thus meaning that their design is likely to be specific to the device dynamics.

In irregular waves, LSPI learns the optimal policy within 6 hours starting from unfavourable conditions, thus proving its superior capacity of learning from a limited set of observations. The same behaviour is observed when the controller is tested in two sea states, alternating every 2 hours. Finally, RL is shown to converge towards a new optimal policy after changing the floater geometry, with the controller still being initialized with the reward function valid for the older system. This proves the adaptive nature of RL control, supporting its ability to account for changes in the system dynamics, e.g. due to marine bio-fouling.

REFERENCES

[1] K. Gunn and C. Stock-Williams, "Quantifying the Potential Global Market for Wave Power," *Proceedings of the 4th International Conference on Ocean Engineering (ICOE 2012)*, pp. 1–7, 2012.

[2] A. F. D. O. Falcão, "Wave energy utilization: A review of the technologies," *Renewable and Sustainable Energy Reviews*, vol. 14, no. 3, pp. 899–918, 2010.

[3] S. H. Salter, J. R. M. Taylor, and N. J. Caldwell, "Power conversion mechanisms for wave energy," *Proceedings of the I MECH E Part M*, vol. 216, no. 1, pp. 1–27, 2002.

[4] J. V. Ringwood, G. Bacelli, and F. Fusco, "Energy-Maximizing Control of Wave-Energy Converters: The Development of Control System Technology to Optimize Their Operation," *IEEE Control Systems Magazine*, vol. 34, no. 5, pp. 30–55, 2014.

[5] A. J. Nambiar, D. I. M. Forehand, M. M. Kramer, R. H. Hansen, and D. M. Ingram, "Effects of hydrodynamic interactions and control within a point absorber array on electrical output," *International Journal of Marine Energy*, vol. 9, pp. 20–40, 2015.

[6] A. Babarit and A. H. Clément, "Optimal latching control of a wave energy device in regular and irregular waves," *Applied Ocean Research*, vol. 28, no. 2, pp. 77–91, 2006.

[7] A. Babarit, M. Guglielmi, and A. H. Clément, "Declutching control of a wave energy converter," *Ocean Engineering*, vol. 36, no. 12–13, pp. 1015–1024, 2009.

[8] F. Fusco and J. V. Ringwood, "A simple and effective real-time controller for wave energy converters," *IEEE Transactions on Sustainable Energy*, vol. 4, no. 1, pp. 21–30, 2013.

[9] T. K. A. Brekken, "On Model Predictive Control for a point absorber Wave Energy Converter," *Proceedings of the IEEE Trondheim PowerTech*, pp. 1–8, 2011.

[10] G. Li and M. R. Belmont, "Model predictive control of sea wave energy converters - Part I: A convex approach for the case of a single device," *Renewable Energy*, vol. 69, pp. 453–463, 2014.

[11] M. Richter, O. Sawodny, M. E. Magaña, and T. K. a. Brekken, "Power optimisation of a point absorber wave energy converter by means of linear model predictive control," *IET Renewable Power Generation*, vol. 8, no. 2, pp. 203–215, 2014.

[12] J. C. Henriques, L. M. Gato, A. F. Falcão, E. Robles, and F. X. Fay, "Latching control of a floating oscillating-water-column wave energy converter," *Renewable Energy*, vol. 90, pp. 229–241, 2016.

[13] F. Fusco and J. V. Ringwood, "Hierarchical robust control of oscillating wave energy converters with uncertain dynamics," *IEEE Transactions on Sustainable Energy*, vol. 5, no. 3, pp. 958–966, 2014.

[14] M. P. Schoen, J. Hals, and T. Moan, "Wave prediction and robust control of heaving wave energy devices for irregular waves," *IEEE Transactions on Energy Conversion*, vol. 26, no. 2, pp. 627–638, 2011.

[15] E. Anderlini, D. I. M. Forehand, P. Stansell, Q. Xiao, and M. Abusara, "Control of a Point Absorber using Reinforcement Learning," *Transactions on Sustainable Energy*, vol. 7, no. 4, pp. 1681–1690, 2016.

[16] M. Eriksson, R. Waters, O. Svensson, J. Isberg, and M. Leijon, "Wave power absorption: Experiments in open sea and simulation," *Journal of Applied Physics*, vol. 102, no. 8, 2007.

[17] M. G. Lagoudakis and R. Parr, "Least-squares policy iteration," *The Journal of Machine Learning Research*, vol. 4, pp. 1107–1149, 2003.

[18] A. Geramifard, T. J. Walsh, S. Tellex, G. Chowdhary, N. Roy, and J. P. How, "A Tutorial on Linear Function Approximators for Dynamic Programming and Reinforcement Learning," *Foundations and Trends® in Machine Learning*, vol. 6, no. 4, pp. 375–451, 2013.

[19] O. Danielsson, "Wave Energy Conversion: Linear Synchronous Permanent Magnet Generator," PhD, Uppsala University, 2006.

[20] M. Eriksson, "Modelling and Experimental Verification of Direct Drive Wave Energy Conversion," PhD, Uppsala University, 2007.

[21] R. Waters, "Energy from Ocean Waves," PhD, Uppsala University, 2008.

[22] M. Stalberg, R. Waters, O. Danielsson, and M. Leijon, "Influence of Generator Damping on Peak Power and Variance of Power for a Direct Drive Wave Energy Converter," *Journal of Offshore Mechanics and Arctic Engineering*, vol. 130, no. 3, pp. 1–4, 2008.

[23] E. Lejerskog, C. Boström, L. Hai, R. Waters, and M. Leijon, "Experimental results on power absorption from a wave energy converter at the Lysekil wave energy research site," *Renewable Energy*, vol. 77, pp. 9–14, 2015.

[24] M. Eriksson, J. Isberg, and M. Leijon, "Theory and experiment on an elastically moored cylindrical buoy," *IEEE Journal of Oceanic Engineering*, vol. 31, no. 4, pp. 959–963, 2006.

[25] D. Forehand, A. E. Kiprakis, A. Nambiar, and R. Wallace, "A Bi-directional Wave-to-Wire Model of an Array of Wave Energy Converters," *IEEE Transactions on Sustainable Energy*, vol. 7, no. 1, pp. 118–128, 2016.

[26] J. Falnes, *Ocean waves and Oscillating systems*, paperback ed. Cambridge University Press, 2005.

[27] G. F. Franklin, J. D. Powell, and A. Emami-Naeini, *Feedback Control of Dynamic Systems*, 6th ed. Pearson, 2008.

[28] I. H. Hutchinson, *A Student's Guide to Numerical Methods*. Cambridge University Press, 2015.

[29] R. S. Sutton and A. G. Barto, *Reinforcement Learning*, hardcover ed. MIT Press, 1998.

[30] L. H. Holthuijsen, *Waves in Oceanic and Coastal Waters*. Cambridge University Press, 2007.

[31] R. Waters, J. Engström, J. Isberg, and M. Leijon, "Wave climate off the Swedish west coast," *Renewable Energy*, vol. 34, no. 6, pp. 1600–1606, 2009.

[32] F. Lewis, D. Vrabie, and K. Vamvoudakis, "Reinforcement Learning and Feedback Control: Using Natural Decision Methods to Design Optimal Adaptive Controllers," *IEEE Control Systems*, vol. 32, no. 6, pp. 76–105, 2012.

Enrico Anderlini received a M.Eng. degree in naval architecture from the University of Southampton, UK, in 2013. Currently, he is working towards the Eng.D. degree in offshore renewable energy at IDCORE, a partnership of the Universities of Edinburgh, Exeter and Strathclyde, UK.

His research interests include marine hydrodynamics, control of wave energy converters, and machine learning.

David I. M. Forehand received the B.Sc. (Hons.) degree in mathematics, in 1990 and the Ph.D. degree in numerical modelling of free-surface waves from the University of Edinburgh, U.K., in 1999 and the M.Sc. degree in applied mathematics from the University of Oxford, U.K., in 1993.

He is a Research Associate with the Institute for Energy Systems, University of Edinburgh, UK. His research interests include non-linear engineering dynamics and the numerical hydrodynamic modelling of wave energy converters and floating bodies.

Elva Bannon Elva Bannon completed a M.Eng. in Advanced Engineering following a B.Eng. (Hons.) in Mechatronic Engineering from Dublin City University.

She then worked in WEC technology research and development focussing on tank testing, modelling and simulation before taking up a post as Senior Research Engineer for Wave Energy Scotland.

Mohammad Abusara received his B.Eng. degree from Birzeit University, Palestine, in 2000 and his Ph.D. degree from the University of Southampton, UK, in 2004, both in Electrical Engineering.

He is currently a Senior Lecturer in Renewable Energy at the University of Exeter, UK. He has over ten years of industrial experience with Bowman Power Group, Southampton, UK, in the field of research and development of digital control of power electronics. During his years in industry, he led the development of a number of commercial products that include grid and parallel connected inverters, microgrid, DC/DC converters for hybrid vehicles, and sensorless drives for high speed permanent magnet machines.