



# Active Flow Control Using Deep Reinforcement Learning with Time-Delays in Markov Decision Process and Autoregressive Policy

DOI:  
[10.1063/5.0086871](https://doi.org/10.1063/5.0086871)

**Document Version**  
Accepted author manuscript

[Link to publication record in Manchester Research Explorer](#)

**Citation for published version (APA):**

Mao, Y., Zhong, S., & Yin, H. (2022). Active Flow Control Using Deep Reinforcement Learning with Time-Delays in Markov Decision Process and Autoregressive Policy. *Physics of Fluids*. <https://doi.org/10.1063/5.0086871>

**Published in:**  
Physics of Fluids

**Citing this paper**

Please note that where the full-text provided on Manchester Research Explorer is the Author Accepted Manuscript or Proof version this may differ from the final Published version. If citing, it is advised that you check and use the publisher's definitive version.

**General rights**

Copyright and moral rights for the publications made accessible in the Research Explorer are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

**Takedown policy**

If you believe that this document breaches copyright please refer to the University of Manchester's Takedown Procedures [<http://man.ac.uk/04Y6Bo>] or contact [uml.scholarlycommunications@manchester.ac.uk](mailto:uml.scholarlycommunications@manchester.ac.uk) providing relevant details, so we can investigate your claim.



# Active Flow Control Using Deep Reinforcement Learning with Time-Delays in Markov Decision Process and Autoregressive Policy

Yiqian Mao (毛逸谦)<sup>1,a)</sup>, Shan Zhong<sup>1</sup>, Hujun Yin<sup>2</sup>

<sup>1</sup>Department of Mechanical, Aerospace and Civil Engineering,

<sup>2</sup>Department of Electrical and Electronic Engineering,

The University of Manchester, Oxford Road, Manchester, M13 9PL, UK

<sup>a)</sup> Author to whom correspondence should be addressed: yiqian.mao@postgrad.manchester.ac.uk

---

## Abstract

Classical active flow control (AFC) methods based on solving the Navier-Stokes equations are laborious and computationally intensive even with the use of reduced-order models. Data-driven methods offer a promising alternative for AFC and they have been applied successfully to reduce the drag of two-dimensional bluff bodies, such as a circular cylinder, using deep reinforcement learning (DRL) paradigms. However, due to onset of weak turbulence in the wake the standard DRL method tends to result in large fluctuations in the unsteady forces acting on the cylinder as the Reynolds number increases. In this study, a Markov decision process (MDP) with time delays is introduced to model and quantify the action delays on the environment in a DRL process due to the time difference between control actuation and flow response along with the use of a first-order autoregressive policy (ARP). This hybrid DRL method is applied to control the vortex shedding process from a two-dimensional circular cylinder using four synthetic jet actuators at a freestream Reynolds number of 400. This method has yielded a stable and coherent control which results in a steadier and more elongated vortex formation zone behind the cylinder hence a much weaker vortex shedding process and less fluctuating lift and drag forces. Compared to the standard DRL method, this method utilizes the historical samples without additional sampling in training and it is capable of reducing the magnitude of drag and lift fluctuations by approximately 90% while achieving a similar level of drag reduction in the deterministic control at the same actuation frequency. This study demonstrates the necessity of including a physics-informed delay and regressive nature in the MDP and the benefits of introducing ARPs to achieve a robust and temporal-coherent control of unsteady forces in active flow control.

**Key words:** Flow control, reinforcement learning, artificial neural networks, autoregressive model, circular cylinder, drag reduction.

---

## 1. Introduction

Flow control has attracted a great deal of research attentions since 1990's because of its potential in improving the aerodynamic performance of transport vehicles beyond a level that could be achieved through optimising the vehicles' geometry alone<sup>1-3</sup>. Passive flow control (PFC) and open-loop active flow control (AFC) methods have been extensively investigated due to their simplicity in implementation and low maintenance need<sup>4-6</sup>. However, these methods cannot always yield desired control effects, especially at off-design conditions<sup>7</sup>. Therefore, more research attentions have begun to be directed towards developing closed-loop AFC methods in recent years in pursuit of higher flow control effectiveness across a larger operation envelope<sup>8-10</sup>.

Due to the high dimensional and nonlinear nature of the governing equations of complex fluid flows, reduced-order models are often required to implement optimization or control with classical AFC methods. Although some physics-based and mathematics-based models<sup>11,12</sup> have shown good feasibility and reliability in practice, achieving a good balance between high fidelity and efficiency remains to be a laborious task. By contrast, machine learning (ML) methods are data-driven and allow engineers to perform AFC without complete prior knowledge of fluid physics. In complex AFC problems with multi-input and multi-output, ML can be applied to develop a computationally efficient surrogate model for predicting control parameters to minimize cost function<sup>13</sup>.

A number of studies have been undertaken to develop effective and reliable ML methods for AFC in recent years, e.g. Shimomura et al.<sup>14</sup>, Li et al.<sup>15</sup> and Fukami et al.<sup>16</sup>. Among these ML methods, supervised learning (SL) and reinforcement learning (RL) have been widely adopted<sup>17-19</sup>. SL aims to establish an optimal model based on existing knowledge and hence requires a sufficient amount of representative labeled data. Recent SL applications in AFC have made use of either artificial neural networks (ANNs), genetic algorithms<sup>20,21</sup> (GAs) or Gaussian process<sup>22,23</sup> (GP) to develop algorithms for modelling complex patterns, flow prediction and control problems. The fully-connected neural networks (FCNNs) and convolutional neural networks (CNNs) are widely-used in ANNs-based AFC methods<sup>24-26</sup>. A comparison of FCNNs and CNNs have been made by Han and Huang<sup>27,28</sup> in opposition control of turbulent channel flows. Some hybrid models combining CNNs with multi-layer perceptron<sup>16</sup> or autoencoder (AE)<sup>29-31</sup> have been proposed for spatial reconstruction, dimensionality reduction and flow estimation in the flow field. To describe the spatial-temporal flow evolution, a combined model of long short-term memory (LSTM) networks and CNN-AE has been proposed<sup>32</sup> and tested<sup>33</sup> for its capability and robustness in controlling the flow around a circular cylinder at a Reynolds number ( $Re$ ) from 20 to 160.

In most AFC tasks with SL, controls are hard to be labeled as 'correct' or 'wrong' even by an expert. Besides, the amount of representative labelled data (such as controlled variables and actuator signals) required by SL may

1 be unrealistic to obtain. In contrast, RL does not need the ‘correct’ strategy as supervisory information but  
2 generates its own data by exploring and evaluating actions against a reward function<sup>34</sup>. Benefiting from its capacity  
3 of modelling policies and value functions in complex RL tasks with continuous state and action space, deep  
4 reinforcement learning (DRL) which combines RL and deep learning has been applied to automatically perform  
5 AFC strategies<sup>35–38</sup>. In DRL, an RL agent samples action-state pairs through interacting with an environment and  
6 adopts ANNs as function approximators to estimate a value function or a policy from the sampled histories. Two  
7 popular DRL methods, i.e. proximal policy optimization (PPO)<sup>39</sup> and twin delayed deep deterministic policy  
8 gradients (TD3)<sup>40</sup>, have been widely adopted in AFC to process continuous action space with high dimensions<sup>41–</sup>  
9 <sup>45</sup>.

10 The application of DRL on the AFC of a bluff body is a sequential decision-making process, and the control  
11 difficulty increases with the complexity of the flow state. In view of the PPO method’s satisfactory training  
12 efficiency and little demand for hyperparameter tuning, Rabault et al.<sup>46</sup> applied a DRL paradigm combining a  
13 PPO method with a fully connected ANN to control the Karman vortex street behind a two-dimensional circular  
14 cylinder using two synthetic jets, one located on the top surface and the other on the bottom surface of cylinder,  
15 at a Reynolds number of 100. Directed by a reward function which minimizes the drag while keeping the lift low,  
16 the DRL agent found a control strategy which suppresses the vortex shedding process and achieves a drag  
17 reduction of approximately 8%. The control strategy discovered using the PPO method was capable of  
18 counteracting against the vortex shedding with the synthetic jet actuation hence stabilizing the fluctuations of  
19 aerodynamic forces while reducing drag

20 Adopting learning parallelization through a multi-environment approach proposed by Rabault et al.<sup>47</sup>, Tang  
21 et al.<sup>48</sup> implemented the DRL paradigm to control the vortex shedding process behind a two-dimensional cylinder  
22 using four synthetic jets over a range of Reynolds numbers from 100 to 400. A drag reduction of 20.4% and 33.1%  
23 were achieved at  $Re = 200$  and 400, respectively. However, random fluctuations in drag and lift appeared to  
24 develop at  $Re = 200$  and increase in magnitude at  $Re = 400$ . Tang et al.<sup>48</sup> attributed the unstable control to the  
25 inherent instability of the flow at  $Re=400$ . This is because when  $Re$  exceeds about 260, turbulence starts appearing  
26 in the shear layer and begins to affect in the state space observed in the near wake<sup>49</sup>. They then introduced global  
27 training and smooth action updates to improve the robustness of a DRL paradigm. Although the control becomes  
28 more stable at Reynolds numbers of 200, 300 and 400, a considerable level of irregular fluctuations in drag and  
29 actuation is still present<sup>48</sup>.

1 Ren et al.<sup>50</sup> also encountered a similar problem in applying DRL control at a higher Reynolds number. Using  
2 the lattice-Boltzmann method (LBM) to simulate the flow environment, they applied DRL to the same flow setup  
3 as in Rabault et al.<sup>46</sup> at a Reynolds number of 1000. The DRL agent was able to find effective AFC strategies with  
4 more training required than at  $Re = 100$ . Nevertheless, the temporal variations in drag exhibit much more random  
5 and significant fluctuations at  $Re = 1000$  and this was also attributed by the authors to the presence of the weakly  
6 turbulence in the flow at this higher Reynolds number.

7 Overall, the results from the aforementioned studies have revealed that although a considerable amount of  
8 reduction in the time-averaged drag can be achieved through AFC with DRL, as the Reynolds number increases  
9 the level of temporal fluctuations in drag tends to become increasingly more random and severe. At Reynolds  
10 numbers higher than 200, the standard PPO method tends to find a ‘cheating’ policy under which although the  
11 time-averaged reward is maximized, random fluctuations and sudden large extremes in drag and lift caused by  
12 irregular control appear at some instants. This inferior control policy can be difficult to avoid by modifying the  
13 reward function since an instant return of the reward at each numerical time step is too expensive and not practical.  
14 Due to the appearance of turbulence in the state space, insufficient regression of the ANN of the time series in the  
15 decision process may result in a deteriorating robustness and temporal-coherence of control trained through the  
16 PPO method.

17 In this study, a Markov decision process (MDP) with time delays is introduced to quantify the action delays  
18 in the DRL process due to the time elapse between actuation and response of flow along with the use of a first-  
19 order autoregressive policy (ARP). This hybrid DRL method is applied to control the vortex shedding process  
20 from a two-dimensional circular cylinder using four synthetic jet actuators at a freestream Reynolds number of  
21 400. This method has yielded a stable and coherent control which results in a steadier and more elongated  
22 recirculation zone behind the cylinder hence a much weaker vortex shedding process. Compared to the standard  
23 PPO method, this method utilizes the historical samples without additional sampling in training and it is capable  
24 of reducing the magnitude of drag and lift fluctuations by approximately 90% while achieving a similar level of  
25 drag reduction in the deterministic control at the same actuation frequency. This study demonstrates the necessity  
26 of including a physics-informed time delay in the MDP and the benefits of introducing ARP to achieve a robust  
27 and temporal-coherent control of unsteady forces in active flow control. Reduction of high-level temporal  
28 fluctuations in drag or lift will help to decrease dynamic loads and structural fatigue leading to an improved  
29 structural durability and operational safety<sup>51-53</sup>.

## 2. Methodology and implementation

### 2.1. Numerical setup

The flow configuration in this study is a two-dimensional laminar flow passing around a circular cylinder at a Reynolds number ( $Re$ ) of 400, where  $Re = \bar{U}D/\nu$ ,  $\bar{U}$  is the average velocity of inflow,  $\nu$  is the kinematic viscosity of the fluid, and  $D$  is the cylinder diameter. Tang et al.<sup>48</sup> applied DRL control within the range of  $Re = 100 - 400$  and they obtained the largest amount of drag reduction as well as the highest level of drag force fluctuations at  $Re = 400$ . Therefore,  $Re = 400$  is selected in this study to enable a comparison of flow control results using different DRL strategies. In this study, the length, velocity, time and vorticity are nondimensionalized with  $D$ ,  $\bar{U}$ ,  $t = D/\bar{U}$  and  $1/t$ , respectively.

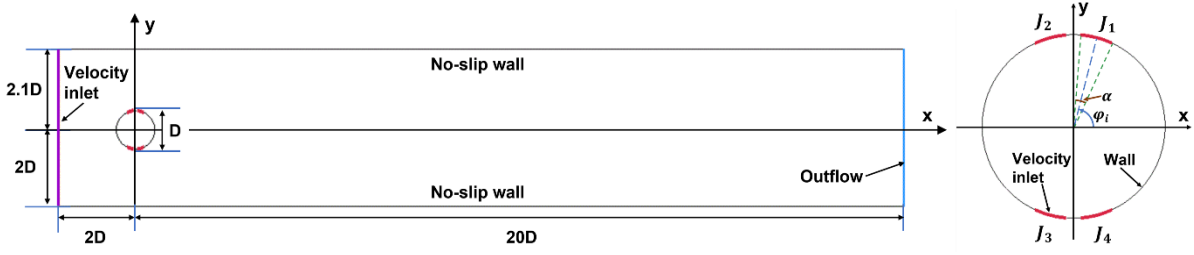
To enable a direct comparison, the same numerical setup used by Tang et al.<sup>48</sup>, which has already been validated for simulating this flow, is also applied in our CFD simulations. The size of the rectangular computational domain is chosen to be  $22D(\text{length}) \times 4.1D(\text{height})$ . As shown in Fig. 1, the cylinder centre is placed with an offset of  $0.05D$  from the centreline of the computational domain to initiate vortex shedding in the simulations at this Reynolds number. A no-slip condition is applied on the top and bottom boundary as well as on the surface of the cylinder. The inflow boundary is located at  $2D$  upstream of the cylinder centre, and a velocity inlet condition with a parabolic velocity profile is applied at this boundary so as to match with the wall boundary conditions imposed on the top and bottom of the computational domain. An outflow boundary condition is set at the outlet boundary.

Four control jets ( $J_1, J_2, J_3, J_4$ ) are symmetrically installed on the cylinder surface with an angular position of  $(\varphi_1, \varphi_2, \varphi_3, \varphi_4) = (75^\circ, 105^\circ, 255^\circ, 285^\circ)$  respectively. The velocity profile at each jet exit follows the cosine distribution defined in Eq. (1) to ensure the velocity continuity between the velocity inlet of jets and no-slip surfaces of cylinder,

$$U_i(\varphi, Q_i) = \frac{\pi}{\alpha D} Q_i \cos\left(\frac{\pi}{\alpha}(\varphi - \varphi_i)\right) \quad (1)$$

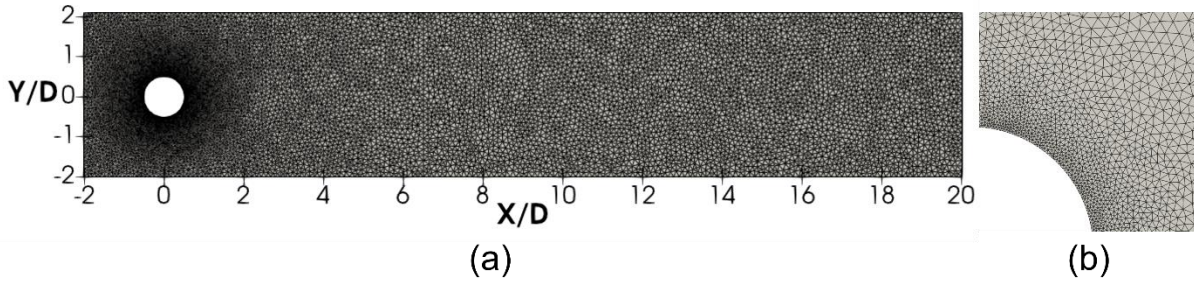
where  $\varphi$  is the angle measured from the positive semi-axis of X-axis,  $\alpha$  is the angular width of each jet which is set as  $\pi/18$ ,  $Q_i$  is the maximum mass flow rate of  $J_i$  taken at its centre location,  $\varphi_i$ . Let  $Q_i^* = Q_i/Q_{ref}$ , where  $Q_{ref}$  is the reference mass flow rate,  $Q_{ref} = \int_{-D/2}^{D/2} \rho U(y) dy$ , in which  $U(y)$  is the inlet velocity profile upstream of the cylinder. A constraint of  $|Q_i^*| \leq 0.05$  is applied to each jet to avoid large power consumption and non-physical actuations. A constraint of zero total mass flow rate of the four jets is applied to avoid adding and

1 subtracting mass directly from the environmental flow. This condition ensures that control strategies learned by  
 2 the agent are more realistic and the numerical scheme is more stable.



3  
 4 **Fig. 1.** Schematics showing computational domain, boundary conditions and control jet locations.

5 In the present work, the flow is assumed to be viscous and incompressible. The incremental pressure  
 6 correction scheme<sup>54</sup> based on the finite element method is adopted to perform the numerical simulations using the  
 7 open-source computing platform FEniCS<sup>55</sup>. The same computational mesh and solution methods used by Tang et  
 8 al.<sup>48</sup> are applied in our study to eliminate discrepancies in results due to the use of different numerical methods.  
 9 An unstructured mesh composed of 25865 triangular grids is established for the domain discretization where the  
 10 grids near the cylinder surface are refined, see Fig. 2. To avoid duplications, the detailed computational setting  
 11 will not be repeated here. For more details of the numerical setup, please refer to Tang et al.<sup>48</sup>



12  
 13 **Fig. 2.** Computation mesh used in the simulations. (a) The whole domain; (b) near the cylinder surface.

14 The drag  $F_D$  and lift  $F_L$  acting on the circular cylinder are obtained via surface integrations of the viscous  
 15 stress and pressure acting around the cylinder. The drag coefficient,  $C_D$ , and the lift coefficient,  $C_L$ , are calculated  
 16 using Eqs. (2) and (3) from  $F_D$  and  $F_L$ . The standard deviation of  $C_D$  and  $C_L$  are defined as  $\sigma_D$  and  $\sigma_L$ ,

$$C_D = \frac{2F_D}{\rho \bar{U}^2 D} \quad (2)$$

$$C_L = \frac{2F_L}{\rho \bar{U}^2 D} \quad (3)$$

1 The Strouhal number,  $St_h = fD/\bar{U}$ , is used to represent the non-dimensional frequency in this study. The  
2 Strouhal number corresponding to the dominant frequency of  $C_D$  and  $C_L$  are denoted as  $St_{h_D}$  and  $St_{h_L}$ . The  
3 Strouhal number corresponding to the period of vortex shedding  $T_S$  is denoted as  $St_{h_S}$ .  
4 The same code for flow simulations used by Tang et al<sup>48</sup> was adopted in this study, for which case validation has  
5 been presented by Tang et al<sup>48</sup>. Nevertheless, we have conducted our own simulations using different mesh  
6 configurations and time-step discretization to validate the spatio-temporal convergence of our simulations and  
7 provided the results in Table 1. The numerical solver adopts the direct solution of lower-upper (LU) decomposition  
8 and hence the iterative residuals do not exist in the computations. Configuration 2 has the identical mesh  
9 configuration and time step used by Tang et al<sup>48</sup>. As it can be seen in the table, the results produced with the  
10 medium mesh (Configuration 2) are within 0.1% of those with the fine mesh (configuration 3) except for a slightly  
11 larger discrepancy in  $St_h$  of less than 0.5%. At to the convergence of time step, when  $\delta t$  is reduced from 0.001 to  
12 0.0005 (the smallest time step) the discrepancies of all results are within 0.1%. Based on our independent study,  
13 the mesh density and time step used in Configuration 2, which are identical to those used by Tang et al, have  
14 produced converged results. Therefore, they are used in the simulations in our work. This will ensure that any  
15 possible discrepancies due to differences in the numerical setup can be eliminated when different DRL methods  
16 are compared.

17 **Table 1** Results of spatiotemporal convergence study.  $\delta t$  is the numerical time step.

Reynold number	Configuration	Mesh resolution	$\delta t$	$\bar{C}_D$	$C_{D,max}$	$ \bar{C}_L $	$C_{L,max}$	$St_{h_S}$	
400	1	Coarse	7166	0.001	3.166	3.501	1.852	2.962	0.3416
	2	Medium	25865	0.001	3.170	3.467	1.810	2.919	0.3417
	3	Fine	128442	0.001	3.171	3.468	1.809	2.919	0.3400
	4	Medium	25865	0.0005	3.168	3.467	1.808	2.918	0.3417
	5	Medium	25865	0.0015	3.172	3.471	1.813	2.925	0.3451

18

## 19 2.2. DRL procedure

20 The DRL framework used in this study is presented in Fig. 3. The DRL methods used in this study are based  
21 on the PPO algorithms implemented through TensorForce API<sup>56</sup>. The DRL framework considers the AFC problem  
22 as a goal-seeking agent interacting with an uncertain flow environment. The agent is a computer program  
23 representing a learner and decision-maker. The environment is the flow system simulated by the CFD solver,  
24 responding to the agent's actions and presenting new states. In each training case, the agent starts observing the  
25 initial state and act when the vortex shedding is fully developed in the uncontrolled wake flow. In each episode,



1 the PPO agent samples a sequence  $\tau$  defined in Eq. (4) composed of  $N_a$  combinations of state ( $s_t$ ), action ( $a_t$ ) and  
 2 reward ( $r_t$ ) through interactions with the environment,

$$\tau = (s_1, a_1, r_1), (s_2, a_2, r_2), \dots (s_{N_a}, a_{N_a}, r_{N_a}), \quad (4)$$

3 where  $N_a$  is the number of times the DRL agent applying the policy  $\pi$  each episode. A small  $N_a$  corresponds to a  
 4 large update interval of action, leading to inefficient sampling and training, while a large  $N_a$  may produce the jet  
 5 actuation too frequently and affect the numerical stability. In the present work,  $N_a$  is set following the experience  
 6 of Tang et al.<sup>48</sup>. The agent pauses learning by  $\tau$  at the end of each episode. Once learning is complete, the agent  
 7 resumes the paused flow environment and proceed to a new episode. The state space consists of the velocity vector  
 8 of 236 probes with the same arrangement as in the study of Tang et al.<sup>48</sup>. The action space consists of four action  
 9 factors corresponding to the mass flow rates of four jets with zero total mass flow rate. The reward is defined by  
 10 the following equation aiming to minimize the drag while keeping the lift low,

$$r = -\langle C_D \rangle_a - w \cdot |\langle C_L \rangle_a| + C \quad (5)$$

11 where  $\langle \cdot \rangle_a$  is the time average over an action time step  $T_a$ .  $w$  is a weighting coefficient used for adjusting the  
 12 pursue of agent in large drag reduction and keeping lift low. In the present work,  $w$  is set as 0.2, the same as in  
 13 the study of Tang et al.<sup>48</sup>.  $C$  is a constant of 4 for plotting more intuitive learning curves in the results and does not  
 14 affect training. Instead of focusing on the current reward, the agent aims to maximize the average cumulative  
 15 reward  $R(t) = \sum_{i>t} \gamma^{i-t} r_i$ , where  $\gamma \in (0,1)$  is a discount factor used for adjusting the interest of the agent to  
 16 focus on long-term or short-term goals and should be set close to but not exceed 1. Following the same setup of  
 17 Tang et al,  $\gamma$  is set as 0.97 in the present study.

18 The PPO algorithms are policy gradient method. The policy function,  $\pi_\theta$  is represented by an ANN where  
 19 all weights are collectively given by the variable,  $\theta$ . Fig. 4 presents the network architectures of the standard and  
 20 hybrid PPO methods. The PPO algorithms in this study have two networks: an actor network and a critic network.  
 21 Both actor and critic networks comprise an input layer, an output layer and two fully connected hidden layers.  
 22 The size of each hidden layer is set to 512, following an empirical test by Rabault et al.<sup>46</sup> with both modeling  
 23 ability and training efficiency.

24 In the actor network, the input is the state set, and the output is the action distribution. The critic network is  
 25 used to approximate the state-value function  $V(s)$ . The state-action value function  $Q(s, a)$  is replaced by an  
 26 advantage function  $A(s, a) = Q(s, a) - V(s)$  to reduce its variability and accelerate training. When the decision-  
 27 making horizon is infinite, the advantage function can be represented as Eq. (6).

$$\hat{A}_t = R(t) - V(s_t) \quad (6)$$

1 where  $\hat{A}_t$  is an estimator of  $A$  at time step  $t$ . The PPO algorithms, proposed by Schulman et al.<sup>39</sup> as a optimization  
 2 of the TRPO algorithms, aim to maximize a “surrogate” objective defined in Eq. (7).

$$L^{CPI}(\theta) = \hat{\mathbb{E}}_t \left[ \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \hat{A}_t \right] = \hat{\mathbb{E}}_t [r_t(\theta) \hat{A}_t] \quad (7)$$

3 where  $\hat{\mathbb{E}}_t[\cdot]$  indicates the empirical expectation over time,  $r_t(\theta)$  denotes the probability ratio of current policy  $\pi_\theta$   
 4 to previous policy  $\pi_{\theta_{\text{old}}}$ . To avoid an excessively large policy update, a clipped surrogate objective defined in Eq.  
 5 (8) is applied to the actor network.

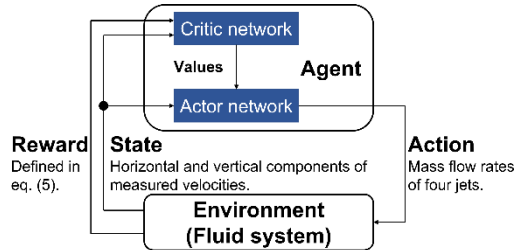
$$L^{CLIP}(\theta) = \hat{\mathbb{E}}_t [\min(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t)] \quad (8)$$

6 where the hyperparameter  $\epsilon$  is chosen as 0.2 as suggested by Schulman et al.<sup>39</sup>. The objective of the critic network  
 7 is defined as  $\hat{\mathbb{E}}_t [-\hat{A}_t^2]$ .

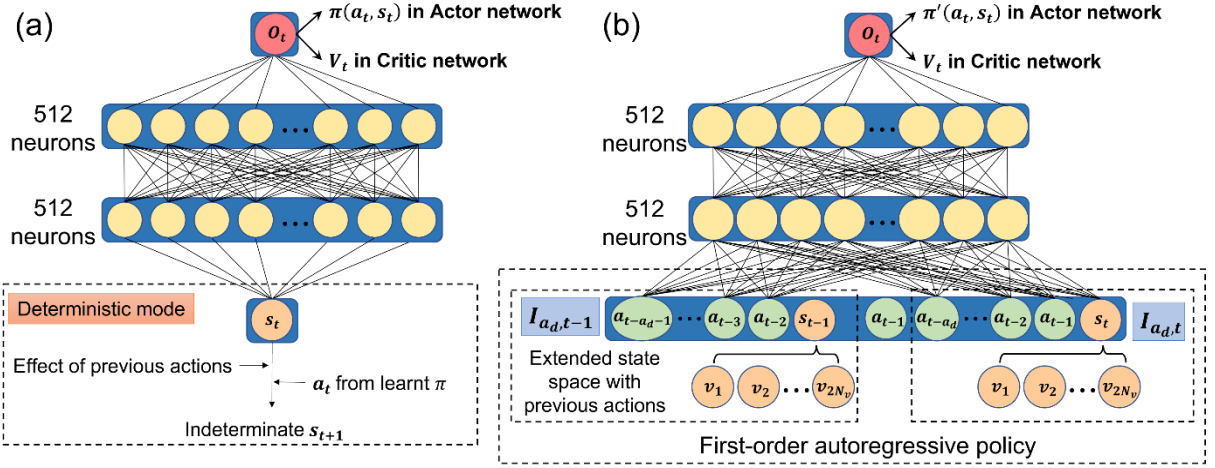
8 In the present study, the duration of each episode is set as  $T_e = 20$ , and the number of action updates per  
 9 episode  $N_a$  is set to 200.  $T_a = T_e/N_a$  equals to 100 numerical time steps, corresponding to 3.3% $T_S$ , which is the  
 10 shedding period in the baseline case. Between two action updates, the value of jet control at each numerical time  
 11 step is given by a linear interpolation defined in Eq. (9) to avoid jumps of computing pressure and velocity in the  
 12 numerical procedure.

$$j_{l,k} = a_{l-1} + k \frac{a_l - a_{l-1}}{N_n} \quad (9)$$

13 where  $N_n = 100$  is the total number of numerical time steps for an action time step.  $k = 1, 2 \dots N_n$  is the current  
 14 numerical time step of the current action.  $l = 1, 2 \dots N_a$  is the current action time step.  $j_{l,k}$  is the control value of  
 15 the jet at  $k$ -th numerical time step of  $l$ -th action time step.  $a_l$  is the action space updated by DRL agent for the  $l$ -  
 16 th action time step. A multi-environment approach proposed by Rabault and Kuhnle<sup>47</sup> is used for parallelizing the  
 17 training into 32 processes.



18  
 19 **Fig. 3.** Schematic of the DRL framework in the present work



**Fig. 4.** Network architectures. (a) Standard PPO method. Since the effect of previous actions is not considered in the modelling of policy  $\pi$ , the environmental state transition may become random due to untreated action delays in the deterministic mode. (b) Proposed hybrid PPO method. A first-order autoregressive policy is applied to the extended state space with previous actions.  $N_p$  is the number of velocity probes.

### 2.3. MDP with delays

In real-world applications of reinforcement learning, three types of delays may exist in the MDPs<sup>57</sup>. They are: (1) observation delay, which exists when the state is not observed immediately; (2) action delay, which exists when an action does not affect the environment immediately but with a time delay; and (3) cost delay, which exists when the action-induced reward is not fully collected till after a certain time has elapsed.

An MDP with constant time delays in observation, action and cost can be denoted as a seven-tuple, i.e.  $M_d = \langle S, A, P_A, r, o_d, a_d, c_d \rangle$ , where  $S$  and  $A$  are the state and action sets, respectively,  $P_A$  is the transition probability, and  $o_d$ ,  $a_d$  and  $c_d$  represent the time delays in observation, action and cost, respectively. In this study, the states and the rewards are collected instantly leading to no observation and cost delays. However, an action delay is inherently present because a jet actuation taking place on the cylinder surface cannot affect the wake flow immediately since a finite period of time is required for its effect to propagate downstream. Therefore, the aforementioned MDP with constant action delays can be reduced to a five-tuple, i.e.  $M_{ad1} = \langle S, A, P_A, r, a_d \rangle$ . Simple ANN is not able to model  $a_d$  in the decision process on its own. As shown in Fig. 4 (a), in the standard PPO method, since the effect of the previous actions is not considered in the modelling of policy, the environmental state transition becomes random for the same current state and action. For example, in the present AFC system, since the flow environment is influenced by the current and previous actions, if an agent under a deterministic policy acts only based on an isolate observation of the current state  $s_t$  at a discrete time step  $t$ , the environment may transit to different  $s_{t+1}$ . Such learnt policy is bound to result in fluctuations in lift and drag,

1 especially when turbulence occurs in the near wake. Therefore, a sufficient ANN's modelling of the control system  
2 should account for these action delays.

3 In the regression processes shown in Fig. 4 (a), the action distribution and the advantage at each time step  
4 are determined by the current states in an MDP without delays. However, when an action delay exists the  
5 information of previous jet actuation should be taken into account in the agent's current decision. Therefore, at  
6 each discrete time step the previous control actions should be added into the input layer of each network to correct  
7 the state-action transition and advantage estimation. According to Katsikopoulos and Engelbrecht<sup>57</sup>,  $M_{ad1}$  can be  
8 converted into an MDP without time delays, i.e.  $M_{ad2} = \langle I_{ad}, A, P_A, r \rangle$  through state space augmentation,  
9 where  $I_{ad} = (a_{t-a_d}, \dots, a_{t-1}, s_t)$  and  $a_{t-i}$  with  $i = 1, 2, 3, \dots, a_d$  are a series of action sets, i.e. the jet  
10 actuations which took place before the present time at which the state,  $s_t$ , is collected.

11 Given enough time a jet actuation will eventually reach the outlet of the computational domain. However, to  
12 ensure the efficiency of the training process the number of previous jet actuations to be included in the MDP has  
13 to be limited. In this study, these previous jet actuations are taken to be those which would have affected the near  
14 field flow within the vortex formation zone. This is a reasonable assumption since the lift and drag, which are  
15 used for defining the reward function, are mainly affected by the wake flow in the vortex formation zone. In the  
16 present work, the vortex formation length  $L_f$  is defined as the streamwise distance from the centre of cylinder to  
17 the point of maximum  $\overline{u'_x u'_x}$ , where  $u'_x$  is the fluctuation in the streamwise component of velocity with respect to  
18 the time averaged flow<sup>58,59</sup>. For the low Reynolds number flow studied in this paper,  $L_f$  of uncontrolled flow is  
19 calculated as 0.98. As reported in previous studies<sup>48,50</sup>, the AFC strategies discovered that using DRL may increase  
20 the vortex formation length. Thus, a longer distance  $d_a = 6L_f$  is selected to avoid missing delay information  
21 during the learning process. To estimate the amount of action delay in time, the convection velocity of jet actuation  
22 is also required and can be reasonably taken as the convection velocity of shedding vortices from the cylinder.  
23 The convection velocity of shedding vortices is found to decrease rapidly in the flow direction within a distance  
24 of  $6D$  and it then becomes nearly constant<sup>60</sup>. To aid the determination of the averaged convection velocity of the  
25 shedding vortices in the vortex formation region, two velocity probes are placed at  $(1, 0)$  and  $(5, 0)$ . Once the lift  
26 of the cylinder becomes periodically stable in the uncontrolled case, the velocity histories containing 20 vortex  
27 shedding cycles are sampled by the two probes. The averaged vortex convection velocity in the vortex formation  
28 region,  $\bar{u}_s$ , can be deduced given the axial spacing between the two probes and the time that it has taken the  
29 vortices to travel from the upstream probe to the downstream one found by a cross-correlation analysis of the  
30 velocity histories. With both the vortex formation zone length and the vortex convection velocity determined, the

1 time of action delay can now be calculated, i.e.  $t_d = (d_a - D/2)/\bar{u}_s$ . Finally, the normalised action delay  $a_d$  can  
 2 be determined as  $a_d = t_d/T_a$ . In the present work,  $a_d = 63$ , and therefore the extended state space can be written  
 3 as  $I_{a_d} = (a_{t-63}, \dots, a_{t-1}, s_t)$ .

#### 4 **2.4. ARPs for continuous control DRL**

5 In continuous control DRL, stochastic policy gradient algorithms generally depend on exploration with  
 6 continuous probability distribution, such as Gaussian distribution, to discover new control strategies. However,  
 7 Gaussian policies do not usually generate samples with temporal coherence thus cannot provide explorations with  
 8 smooth trajectories corresponding to secure and rewarding behaviours in most practical continuous control tasks.  
 9 In addition, due to the lack of historical information and experience memory in the state, the action update of the  
 10 policy may not lead to effective exploration of the environment, and it becomes more and more inefficient as the  
 11 action rate increases. Korenkevych et al.<sup>61</sup> introduced a series of stationary autoregressive processes (ARPs) to  
 12 promote agent’s exploration in continuous RL control tasks. These processes exhibit two main characteristics:  
 13 subsequent process observations are temporally coherent with a continuously adjustable degree of coherence, and  
 14 the process stationary distribution is normal.

15 When the complexity of the AFC problem increases, the sampling inefficiency of DRL and complicated state  
 16 and action spaces may result in a more difficult exploration of the agent. A regression or time series model is  
 17 expected to improve the utilization efficiency of samples and improve the temporal prediction of ANN. Therefore,  
 18 the present work deploys a modified first-order ARP into the DRL paradigm. In the Gaussian policy, the actions  
 19 are sampled through  $\pi_\theta(a_t, s_t) \sim \mathcal{N}(\mu_\theta(s_t), \sigma_\theta^2(s_t))$ .  $\theta$  is the parameter factor of  $\pi_\theta$ .  $\mu_\theta$  and  $\sigma_\theta^2$  are both vectors  
 20 parametrized by ANNs. When the episode is not updated, the policy  $\pi_\theta$  does not change. Sampling by the same  
 21 policy, the action mean will gradually approximate the expectation. We replaced the original observations  $s_t$  with  
 22 extended observations  $\tilde{s}_t = (s_{t-p}, a_{t-p}, \dots, s_{t-1}, a_{t-1}, s_t)$  of an  $p$  order ARP, while the substitute in  
 23 (Korenkevych et al. 2019) is chosen as the white noise component. This modification allows ARP to improve the  
 24 whole action update, rather than the noise term. As illustrated in Fig. 11 (b) later, the modified ARP is aimed to  
 25 find the optimal control with overall temporal coherence. To implement the ARP, a history-dependent policy  
 26  $\pi_\theta(a_t | s_t, h_t^p)$  is defined, where  $h_t^p = (s_{t-p}, a_{t-p}, \dots, s_{t-1}, a_{t-1})$  includes past  $p$  states and actions. The new  
 27 policy form replaces the original form of policy  $\pi_\theta(a_t, s_t)$  during the descent process of policy gradient, which is  
 28 a Markov stochastic process. According to the Markov property, the probability distribution of future states  
 29 depends only on the current state but not on the entire historical path. Therefore, history-dependent policies with  
 30 state space  $h_t^p$  do not necessarily induce Markov stochastic processes even if the environmental transition

1 probabilities are Markovian<sup>62</sup>. However, with an extended state space  $(h_t^p, s_t)$  where  $h_t^p$  is fixed in size, such  
 2 history-dependent policy induces a Markov stochastic process. Denote  $\tilde{M}^p = (\tilde{S}, \tilde{A}, \tilde{P}_A(\cdot | a, \tilde{s}), \tilde{r}(\tilde{s}, a))$  as a  
 3 modified MDP with  $p$  order ARP.  $\tilde{A} = A$ ,  $\tilde{S}$ ,  $\tilde{P}_A(\cdot | a, \tilde{s})$  and  $\tilde{r}(\tilde{s}, a)$  are defined as follows:

$$\forall \tilde{s}, \tilde{s}' \in \tilde{S}:$$

$$\tilde{s} = (s_1, a_1, \dots, s_p, a_p, s_{p+1}), a_k \in A, s_k \in S \forall k \quad (10)$$

$$\tilde{s}' = (s'_1, a'_1, \dots, s'_p, a'_p, s'_{p+1}), a'_k \in A, s'_k \in S \forall k \quad (11)$$

$$\tilde{P}(\tilde{s}' | a, \tilde{s}) = \begin{cases} \tilde{P}(s'_{p+1} | a, s_{p+1}), & \text{if } s'_k = s_{k+1}, k \leq p \\ & a'_k = a_{k+1}, k < p \\ & a'_p = a \\ 0 & \text{otherwise} \end{cases} \quad (12)$$

$$\tilde{r}(\tilde{s}, a) = r(s_{p+1}, a) \quad (13)$$

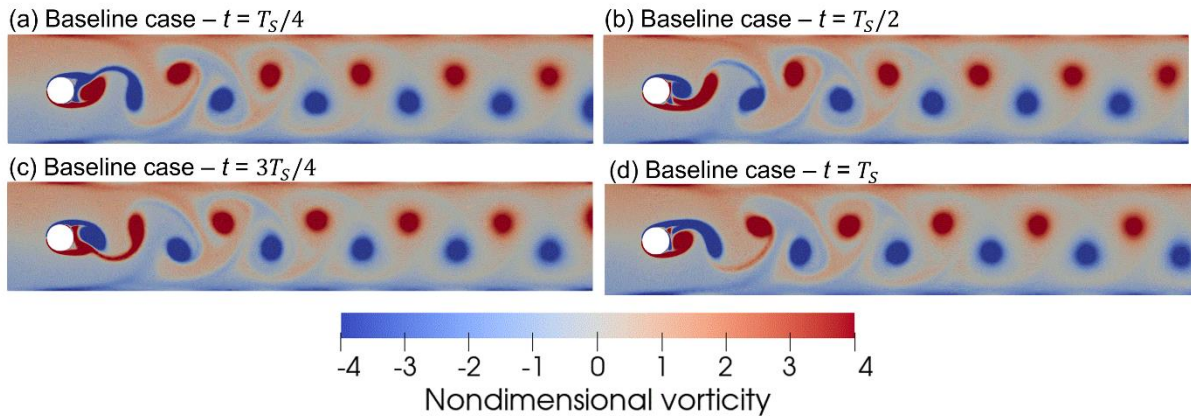
4 The applicability of existing learning algorithms on  $\tilde{M}^p$  has already been validated by Korenkevych et al<sup>61</sup>.  
 5 Therefore, it can be implemented into the present PPO algorithms through augmentation of the state space. In  
 6  $\tilde{M}_{ad2}^p$ , the original state space  $I_{a_d,t}$  in  $M_{ad2}$  is then replaced by  $\tilde{I}_{a_d,t} = (I_{a_d,t-p}, a_{t-p}, \dots, I_{a_d,t-1}, a_{t-1}, I_{a_d,t})$ . A  
 7 corresponding initial state of  $\tilde{M}_{ad2}^p$  is defined as  $\tilde{I}_{a_d,0} = (I_{a_d,0}, a_0, \dots, I_{a_d,0}, a_0, I_{a_d,0})$ , where  $I_{a_d,0}$  is the initial  
 8 states of  $M_{ad2}$ ,  $a_0$  is any element of action set  $A$  because it does not affect future transitions and rewards. The  
 9 sequence  $\tilde{I}_{a_d,t}$  considers time-delays and possesses better temporal coherence in the regression process of NNs  
 10 but does not change the noise proportion in actions. In the present work, we use a first-order ARP in which  $p = 1$   
 11 and  $\tilde{I}_{a_d,1} = (I_{a_d,t-1}, a_{t-1}, I_{a_d,t})$  as shown in Fig. 4 (b).

## 12 3. Results and Discussion

### 13 3.1. The baseline flow

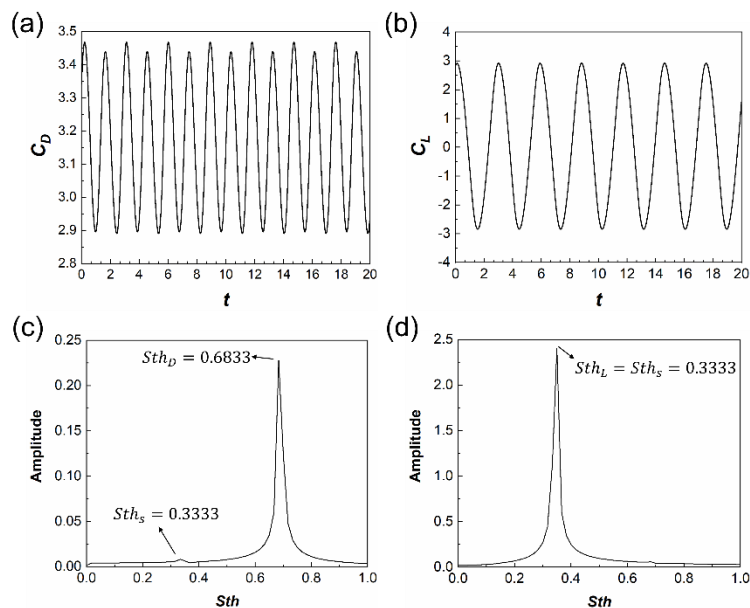
14 Before the evaluation of the performance of different DRL methods in training agents and implementing  
 15 active control, a baseline case of uncontrolled flow is performed first. The flow passing around the cylinder  
 16 without active control (i.e.,  $Q_i^* = 0$ ,  $i = 1,2,3,4$ ) is simulated, and the baseline case is monitored for 20  
 17 nondimensional time when  $C_D$  and  $C_L$  begin to vary periodically. Figure 5 (Multimedia view) presents the  
 18 instantaneous and time averaged vorticity contours and a video showing the flow evolution in the baseline case.  
 19 Figure 5 (a) is captured at the time corresponding to the minimum  $C_L$  during a vortex shedding period,  $T_S$ , and the  
 20 subsequent snapshots are captured at every quarter of a cycle. All results of velocity and vorticity are normalized

1 by  $\bar{U}$  and  $\bar{U}/D$  respectively, in this study. In a vortex shedding period, two vortices with an opposite rotation are  
 2 shed in an alternate manner, leading to periodic pressure variations on the cylinder surface.



3  
 4 **Fig. 5.** Contours of instantaneous vorticity at  $t =$  (a)  $T_S/4$ , (b)  $T_S/2$ , (c)  $3T_S/4$  and (d)  $T_S$  in the baseline case.  
 5 (Multimedia view)

6 Figure 6 illustrates the time series of  $C_D$  and  $C_L$  in the baseline case and their frequency spectra obtained with  
 7 Fast Fourier Transform (FFT). The time averaged  $C_D$  and  $C_L$  over a period are 3.17 and 0.0273, respectively. Note,  
 8 to enable a direct comparison of the oscillatory properties of the variables of interest, the time averaged value of  
 9  $C_D$  and  $C_L$  have been subtracted before the FFT analysis. The vortex shedding frequency corresponds to  $St h_D =$   
 10 0.333. In a single vortex shedding period,  $C_L$  experiences one cycle, while  $C_D$  experiences two cycles, i.e.  $St h_D =$   
 11  $0.683 \approx 2St h_S$ . Due to a slight difference in the vortices shed from the top and bottom of the cylinder, the two  
 12 circles of  $C_D$  becomes different, resulting in a local maximum obtained at  $St h = St h_S$  in the spectrum of  $C_D$ .  
 13 Furthermore,  $C_L$  is asymmetrical about the time axis.



1 **Fig. 6.** Temporal variation of  $C_D$  (a) and  $C_L$  (b) and FFT analyses of  $C_D$  (c) and  $C_L$  (d) in the baseline case. In

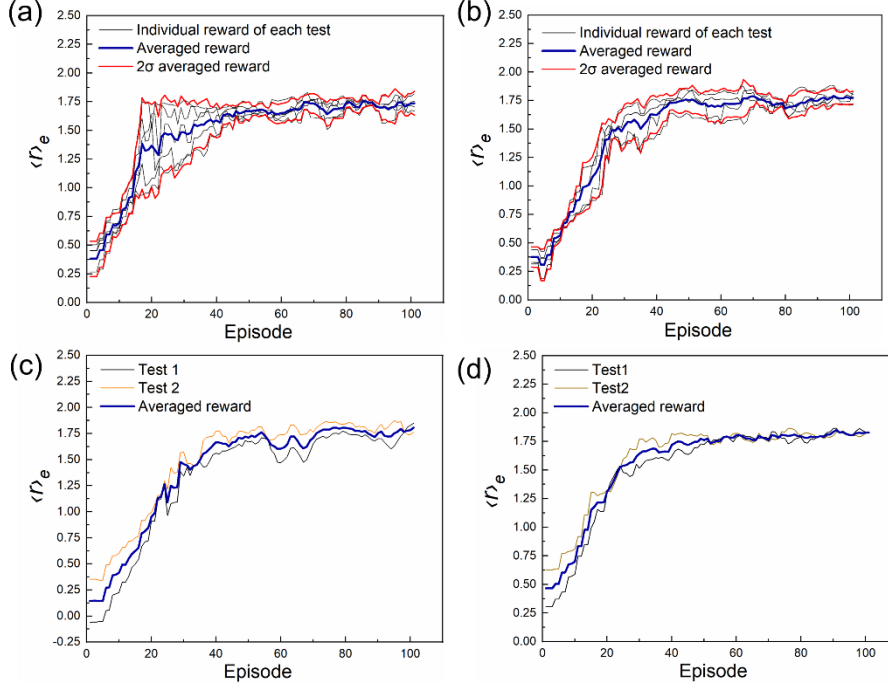
2 FFT, the time series of  $C_D$  and  $C_L$  have their temporal mean subtracted.

### 3 **3.2. Responses of $C_D$ and $C_L$ to different DRL methods**

4 In the present work, four methods are tested and compared for their performance in AFC. The test consists  
5 of two phases: training and deterministic control. In the training phase, each method is used to train an agent until  
6 the policy learned by the neural networks of the agent reaches the convergence that is the reward has nearly  
7 stopped increasing as the episode number increases. After meeting the convergence, all agents have been trained  
8 to episode 3200. To accelerate training, each agent simultaneously samples the state-action-reward sets in 32  
9 parallel environments to update a single policy. After training, the last environment is selected for comparison.  
10 For the four DRL methods tested, we denote the vanilla PPO algorithms as the PPO method, the PPO method  
11 considering delays in the MDP as the DMDP-PPO method, the PPO method employing first-order ARP as the  
12 ARP-PPO method, and the PPO algorithms combining MDP with delays and first-order ARP as the ARP-DMDP-  
13 PPO method. To avoid random effect, the PPO and ARP-DMDP-PPO methods have been used to perform training  
14 five times. While both ARP-PPO and DMDP-PPO methods have been adopted for training two times, due to their  
15 high training costs. Averaged performances are plotted for comparison.

16 In training, the average reward per episode is defined as  $\langle r \rangle_e$ . The averaged learning curve represented by  
17 variations of  $\langle r \rangle_e$  with episode number and the confidence interval are illustrated in Fig. 7. All four methods can  
18 be applied to achieve robust training. In the initial stage of training, the average rewards of PPO and ARP-DMDP-  
19 PPO increase rapidly with the number of episodes. During the main exploration stage, the learning curve of PPO  
20 has a large confidence interval corresponding to the training instability exceeding the effect of exploration. The  
21 convergence differences by more than 20 episodes may indicate that the learnt policies finally reaching the  
22 maximum reward in different training repeats are inconsistent, which may also reflect the insufficient ability of  
23 ANN in time series regression for approximating the policy. The learning curves indicate that learning occurs  
24 consistently in about 50 episodes, then keeps tuned to 100 episodes to reach fine convergence. The maximum  $\langle r \rangle_e$   
25 of the averaged learning curves through the PPO and ARP-DMDP-PPO method is 1.75 and 1.79, respectively.  
26 Since the ARP-DMDP-PPO method improves the ANN's regression through embedding the historical data  
27 without additional sampling, it requires only a slightly larger training time than the PPO method under the same  
28 number of episodes.





**Fig. 7.** Learning curves of different methods in the last environment to indicate the training robustness. For each DRL method, the training repeats have been performed using the same hyperparameters but different random seeds. For the PPO (a) and the ARP-DMDP-PPO method (b) with five training repeats, the averaged learning curve and the 95% confidence interval are presented. For the ARP-PPO (c) and DMDP-PPO method (d) with two training repeats, the averaged learning curve is presented.

In the deterministic control phase, each agent selects the best policy (the policy that achieved the greatest reward) in its training history to start control from the same initial state of the baseline case and updates the jet actuation with the same frequency as action updates in training. All results in the control phase are obtained by agents deterministically acting without exploration. The actions are not sampled for training since all agents do not update their policy during the control phase. The nondimensional time  $t$  of a complete control process is 120. Figure 8 and 9 presents the temporal variations of  $C_D$  and  $C_L$  controlled by four agents, respectively. It is noted that all agents produce a significant decrease in  $C_D$  within  $t = 20$ . The PPO and ARP-PPO agents cannot maintain a stationary  $C_D$  since their policies are learned by the networks without information of action delays. By contrast, when considering action delays,  $C_D$  controlled by the DMDP-PPO and ARP-DMDP-PPO agents has become nearly stationary after  $t = 40$ . Table 2 summarizes  $\sigma_D$  and  $\sigma_L$  controlled by four agents after  $t = 40$ .

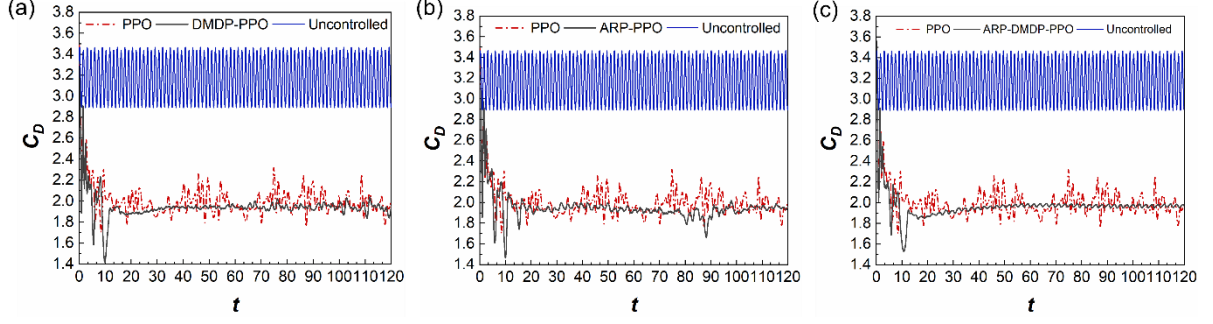
The drag reduction in control is defined as  $(\langle C_D \rangle_b - \langle C_D \rangle_d) / \langle C_D \rangle_b$ , where  $\langle C_D \rangle_d$  is the time averaged  $C_D$  in the DRL cases. Similar to the results of the same flow configuration by Tang et al<sup>48</sup>, the PPO agent has obtained a rapid  $C_D$  reduction, but then  $C_D$  and  $C_L$  keep fluctuating significantly and randomly in the remaining time. The

1 drag reduction by the PPO agent is 36.8%, greater than the drag reduction of 33.1% achieved by Tang et al<sup>48</sup>. This  
2 difference of control may indicate the training of the PPO agent in Tang et al<sup>48</sup> can still be tuned by more episodes  
3 to reach a better convergence. As the number of episodes increases, the best policy during training will be more  
4 probable to close to the optimal policy. Therefore, we choose a higher number of training episodes than Tang et  
5 al. based on the available computational resources to avoid the potential influence of insufficient training. The  
6 result shows that sufficient training of PPO method still cannot find a stable control strategy to reduce the  
7 fluctuations of  $C_D$  and  $C_L$  in subsection 3.3.

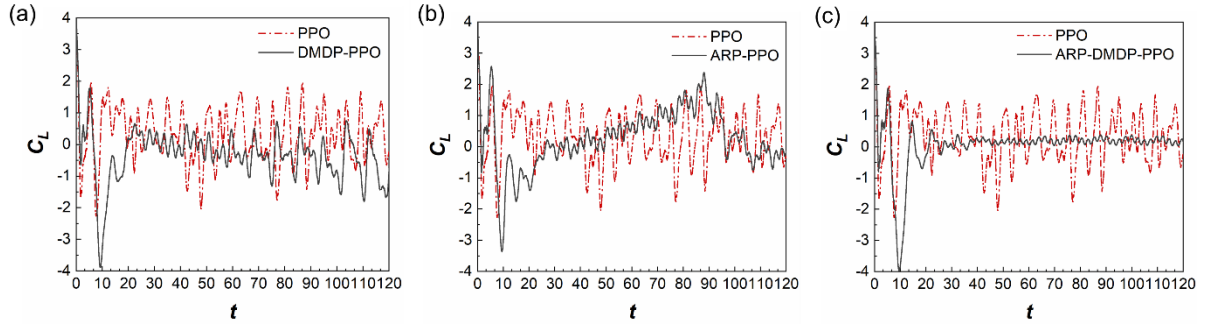
8 After considering action delays in MDP, the DMDP-PPO agent controls  $C_D$  and  $|C_L|$  rapidly decrease, and  
9 then fluctuate with a much lower  $\sigma_D$  and  $\sigma_L$  than the PPO agent. After the rapid decrease of  $C_D$ , the DMDP-PPO  
10 agent tends to find a stationary control with small fluctuations at the beginning. But the opposite is true, and  $C_D$   
11 and  $C_L$  fluctuate in an increasing variance, which indicates the control effects cannot be maintained but gradually  
12 worsen as fluctuation accumulate.

13 Historical states and actions are valuable for the neural networks to estimate advantages  $A_t$  and output  
14 smooth and temporal coherent actions. Autoregressive models are popular and effective tools to perform a  
15 regression of the variable utilizing past values. Here, the ARP-PPO method and PPO method are compared to  
16 investigate the effects of first-order ARP. As shown in Fig. 8 (b) and Fig. 9 (b), after the rapid decrease in the  
17 initial time,  $C_D$  and  $C_L$  controlled by the ARP-PPO agent fluctuate in a lower variance than the PPO agent, but  
18 they are non-stationary with an unsteady trend. As shown in Table 2, the ARP-PPO agent produces smaller  $\sigma_D$   
19 and  $\sigma_L$  than the PPO agent but larger  $\sigma_D$  and  $\sigma_L$  than the DMDP-PPO agent. The performance of the ARP-PPO  
20 agent illustrates the positive effect of ARP, but also indicates that ARP cannot replace DMDP to quantify delays.  
21 Therefore, ARPs should be deployed into the DMDP-PPO method instead of directly adopted in the PPO method.

22 Finally, a hybrid ARP-DMDP-PPO method is used to implement effective control in the AFC system of this  
23 study. After the ARP-DMDP-PPO agent performs control,  $C_D$  rapidly decreases to reach convergence and begins  
24 to change with a small amplitude after  $t = 40$ . The drag reduction by the ARP-DMDP-PPO agent is 38%, which  
25 is higher than the drag reduction obtained by the PPO agent. The ARP-DMDP-PPO agent spends the same time  
26 as the PPO agent to reach around the minimum  $C_D$  but maintains a steadier flow state with a much lower  $\sigma_D$  and  
27  $\sigma_L$  of 87% and 90%, respectively.



**Fig. 8.** Real-time control by different agents showing temporal variations in  $C_D$ . (a) the DMDP-PPO agent, (b) the ARP-PPO agent, (c) the ARP-DMDP-PPO agent.



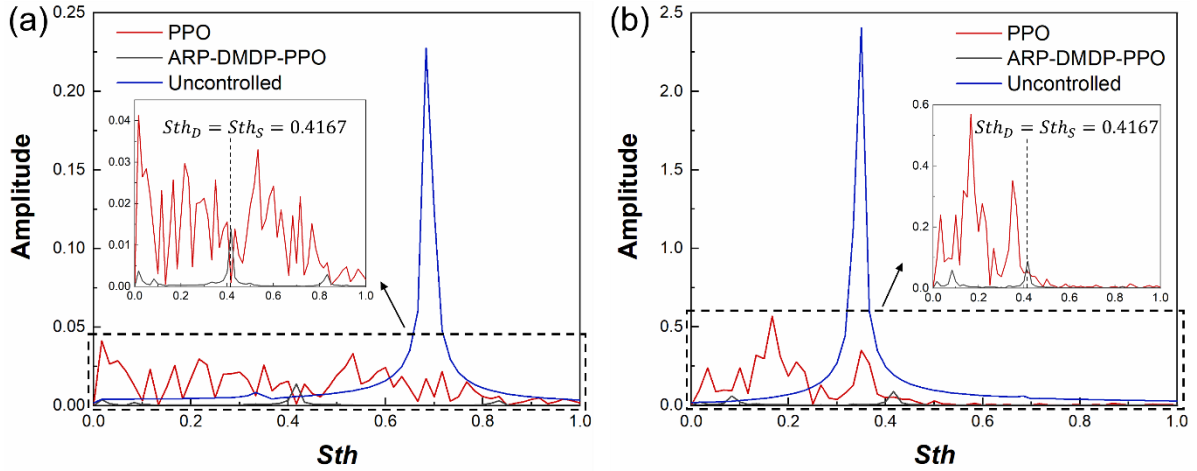
**Fig. 9.** Real-time control by different agents showing temporal variations in  $C_L$ . (a) the DMDP-PPO agent, (b) the ARP-PPO agent, (c) the ARP-DMDP-PPO agent.

**Table 2** Comparison of  $\sigma_D$  and  $\sigma_L$  sampled from  $t = 40$  to  $120$  in four controlled cases.

Agent	$\sigma_D$	$\sigma_L$
PPO	0.0928	0.811
DMDP-PPO	0.0273	0.546
ARP-PPO	0.0464	0.663
ARP-DMDP-PPO	0.0121	0.081

Figure 10 presents the FFT of the time series data of  $C_D$  and  $C_L$  sampled from  $t = 60$  to  $120$  in the PPO case and ARP-DMDP-PPO case.  $C_D$  and  $C_L$  are subtracted by their mean value before FFT analysis. Similar to the results of Tang et al.<sup>48</sup>, amplitudes of  $C_D$  and  $C_L$  at their dominant frequencies have been significantly reduced by the PPO agent. However, the PPO agent has introduced a considerable amount of fluctuation of  $C_D$  and  $C_L$  at the Strouhal numbers lower than  $St_{h_D}$  and  $St_{h_L}$  in the baseline case. Through the modification of jet actuation

1 controlled by the ARP-DMDP-PPO agent,  $St_{h_D}$  and  $St_{h_L}$  become consistent with  $St_{h_1}$  corresponding to the  
 2 actuation frequency presented in Fig. 11. Compared to the baseline case, the ARP-DMDP-PPO agent has led to  
 3 not only the dominant frequency being shifted to a slightly higher frequency but also its amplitude being reduced  
 4 substantially. Compared to the PPO agent, the ARP-DMDP-PPO agent suppresses the amplitude of  $C_D$  and  $C_L$   
 5 more effectively and meanwhile avoids introducing unnecessary random interference at other Strouhal numbers.



6

7 **Fig. 10.** FFT analysis of  $C_D$  and  $C_L$  sampled from  $t = 60$  to  $120$  in the baseline case and two controlled cases.

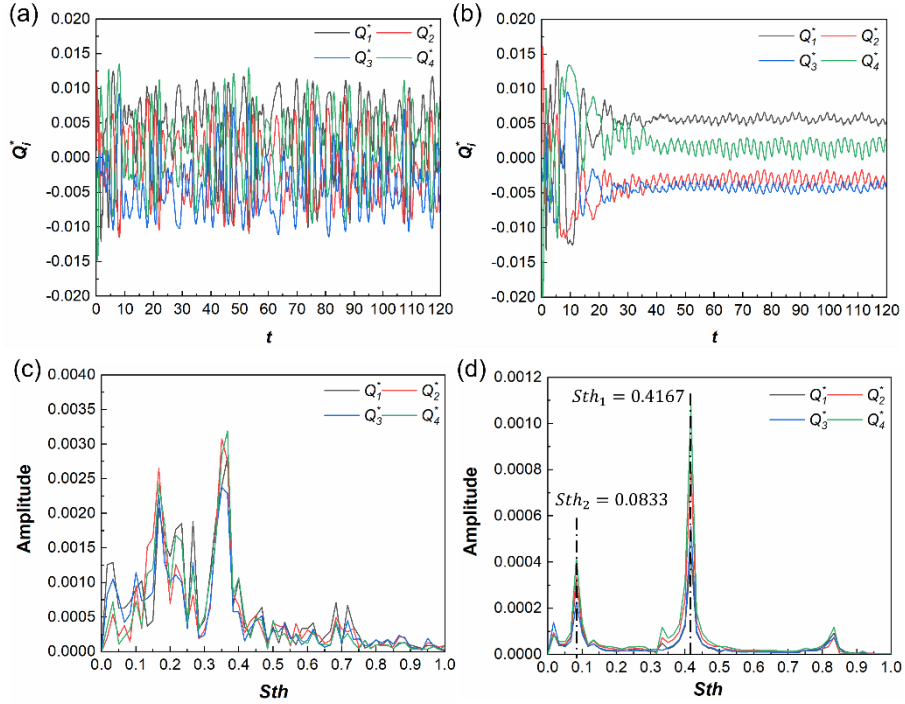
8

Time series of  $C_D$  and  $C_L$  have been subtracted by their temporal mean. (a)  $C_D$ . (b)  $C_L$ .

9

A realistic AFC system usually requires DRL agents to output smooth and temporal coherent jet control due  
 10 to the physical constraints of jet actuators. Figure 11 compares the mass flow rates of all jets and their FFT between  
 11 the PPO case and ARP-DMDP-PPO case during control. The FFT is sampled from the time series data of  $Q_1^*$  to  
 12  $Q_4^*$  (nondimensional mass flow rates of four jets) from  $t = 60$  to  $120$  in both cases. To directly compare the pure  
 13 oscillatory properties of the variables of interest,  $Q_1^*$  to  $Q_4^*$  are subtracted by their mean value before the FFT  
 14 analysis. The PPO agent keeps outputting stochastic jet signals during the entire control process. Similar to the  
 15 control strategy discovered by Rabault et al. at  $Re = 100$ , the ARP-DMDP-PPO agent changes the flow  
 16 configuration into a low-drag state through large-amplitude jet actuation. From  $t = 20$  to  $t = 40$ , since the main  
 17  $C_D$  reduction has already been obtained, each jet signal gradually converges and oscillates to slightly modify the  
 18 ambient flow. After  $t = 40$ , all jet signals begin to settle in a stationary trend, resulting in a small pseudo-periodic  
 19 oscillation of  $C_D$  and  $C_L$  shown in Fig. 8 (c) and Fig. 9 (c). All jet signals by the ARP-DMDP-PPO agent have  
 20 almost the same spectral distribution where the amplitudes are mainly concentrated at  $St_{h_1} = 0.4167$  and  $St_{h_2} =$   
 21  $0.0833$ . On the contrary, the PPO agent has introduced many fluctuations at frequencies lower than  $St_{h_s}$  to each

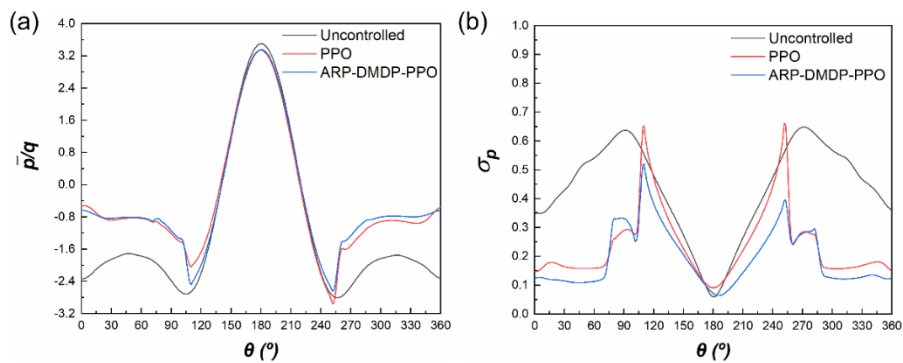
1 jet signal. Furthermore, all control signals are inconsistent with each other and have much larger amplitudes than  
 2 the control signals by the ARP-DMDP-PPO agent.



3  
 4 **Fig. 11.** Temporal variation of  $Q_1^*$  to  $Q_4^*$  (The nondimensional mass flow rates of four jets) in (a) PPO case and  
 5 (b) ARP-DMDP-PPO case. FFT analysis of  $Q_1^*$  to  $Q_4^*$  sampled from  $t = 60$  to  $120$  in (c) PPO case and (d)  
 6 ARP-DMDP-PPO case, the time series have been subtracted by their temporal mean.

7  
 8 Since the pressure drag and lift are the x-integral and y-integral of cylinder surface pressure in the numerical  
 9 computation, the optimization mechanism by the agents is specifically analysed from the temporal variation of  
 10 the pressure distribution on the cylinder surface. Figure 12 presents the time-averaged distribution of pressure and  
 11 the standard deviation of surface pressure on the cylinder surface through the whole control process.  $\theta = 180^\circ$   
 12 corresponds to the stagnation point of the cylinder face.  $\bar{p}$  is the time-averaged pressure normalized by the  
 13 dynamic pressure  $q = \rho \bar{U}^2 / 2$ ,  $\sigma_p$  is the standard deviation of surface pressure. Compared with the baseline case,  
 14 both agents significantly reduce the absolute value of time-average negative pressure and  $\sigma_p$  on the leeward side  
 15 of the cylinder, resulting in a large reduction of  $C_D$  and  $\sigma_D$ . However, the PPO agent cannot stabilize the controlled  
 16 flow because its policy lacks the information of delays and historical states-actions sequences. Compared to the  
 17 baseline case, the large jet actuation in the PPO case increases the pressure fluctuation around  $\theta = 105^\circ$  and  $255^\circ$   
 18 even though it reduces the fluctuation in the flow direction. Furthermore, the PPO agent causes a more asymmetric

1 pressure distribution at the upper and lower surfaces of the cylinder, resulting in a higher time averaged  $C_L$ . On  
 2 the contrary, the ARP-DMDP-PPO agent performs a more precise control to reduce the pressure difference  
 3 between the upper and lower surfaces along the horizontal axis of the cylinder. The ARP-DMDP-PPO agent is  
 4 capable of identifying the asymmetry in the vortex structure in the wake caused by the offset of the cylinder  
 5 relative to the centerline of inflow shown and outputs an asymmetric jet actuation to maintain a symmetric time  
 6 averaged pressure distribution. As a result of the asymmetric jet actuation,  $\sigma_p$  is asymmetrically distributed around  
 7 the cylinder. Compared to the PPO agent, the ARP-DMDP-PPO agent reduces  $\sigma_p$  at most angles especially at  
 8  $\theta = 105^\circ$  and  $255^\circ$ , resulting in a lower fluctuation of drag and lift.

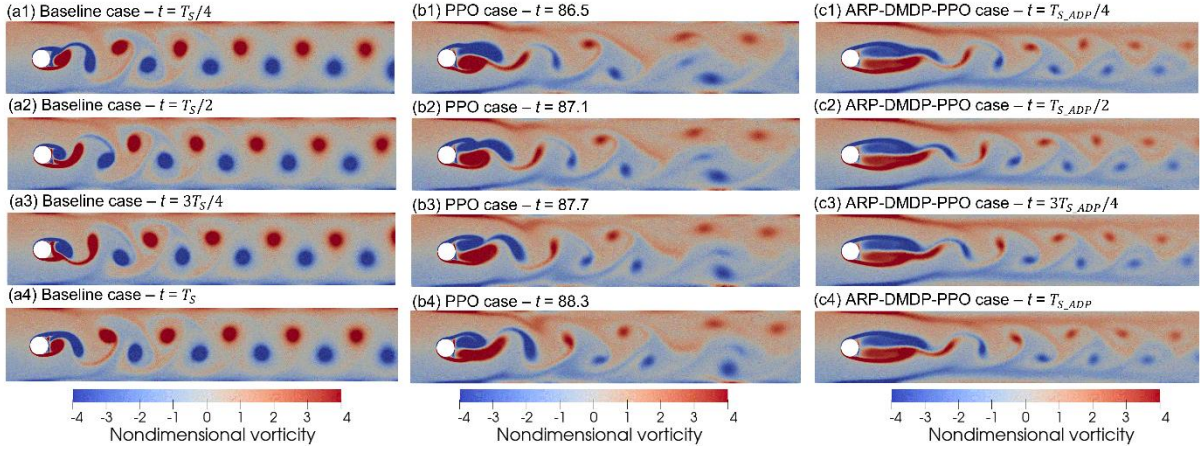


9  
 10 **Fig. 12.** Distribution of  $\bar{p}/q$  and  $\sigma_p$  on the cylinder surface calculated over the whole control process. (a)  $\bar{p}/q$ .  
 11 (b)  $\sigma_p$ .

### 12 3.3. Effects of flow control on vortex formation

13 The effects of controlling jets on the vortex shedding are further analysed in this section. Figure 13  
 14 (Multimedia view) presents the instantaneous vorticity contours and a video of the flow evolution during the entire  
 15 control process in the baseline, PPO and ARP-DMDP-PPO case. In the ARP-DMDP-PPO case, the first image is  
 16 captured at a time corresponding to the minimum  $C_L$  in a single period  $T_{S\_ADP}$ , and the subsequent figures are  
 17 captured every quarter of a cycle. The vorticity contours of the baseline case have been presented in Fig. 13 (a1)  
 18 – (a4). It can be seen clearly that in the baseline case, formation of the shedding vortices takes place within a  
 19 distance much closer to the cylinder than in the controlled cases and the strength of vortices in remains high  
 20 towards the computational domain. In the PPO case, the formation zone has clearly been extended further  
 21 downstream compared to the baseline case and the vortices further downstream are much weaker. The ARP-  
 22 DMDP-PPO agent seems to produce the longest and most steady formation zone with the vortex shedding taking  
 23 place at the end of it. The location of vortex shedding is closer to the centreline of the cylinder and the width of  
 24 the vortex street behind it appears narrower than in the other two cases. A longer and steadier formation zone is

1 expected to lead to a less fluctuating pressure in the near wake resulting in a lower vortex-induced drag. Similar  
 2 phenomena were reported by researchers using passive flow control techniques<sup>63,64</sup>, such as fairings and splitter  
 3 plates.



4

5 **Fig. 13.** Contours of instantaneous vorticity at (a1) - (a4):  $t = T_S/4, T_S/2, 3T_S/4$  and  $T_S$  in the baseline case,  
 6 (b1) - (b4):  $t = 86.5, 87.1, 87.7$  and  $88.3$  in the PPO case and (c1) - (c4):  $t = T_{S\_ADP}/4, T_{S\_ADP}/2, 3T_{S\_ADP}/4$   
 7 and  $T_{S\_ADP}$  in the ARP-DMDP-PPO case. (Multimedia view)

8

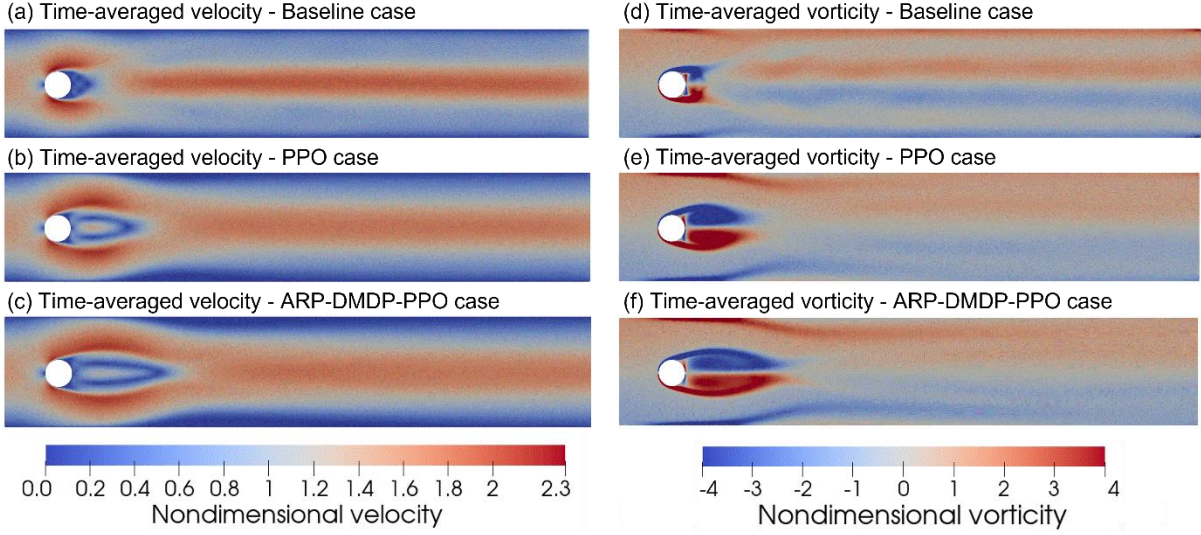
9 Figure 14 presents the contours of time-averaged velocity magnitude and vorticity over the entire control  
 10 process ( $t=120$ ) in the baseline case ( $40T_S$ ), the PPO case and the ARP-DMDP-PPO case ( $50T_{S\_ADP}$ ). In  
 11 comparison to the baseline case, with the jet control the recirculation bubble in the near wake is extended  
 12 substantially in length, resulting in a reduced pressure drop shown in Fig. 12 (a). Compared to the PPO agent, the  
 13 ARP-DMDP-PPO agent gives rise to a more streamlined recirculation zone within which the magnitude of  
 14 vorticity is smaller and the regions of high vorticity are located further away from the cylinder, leading to a steady  
 15 pressure distribution on the cylinder surface.

16

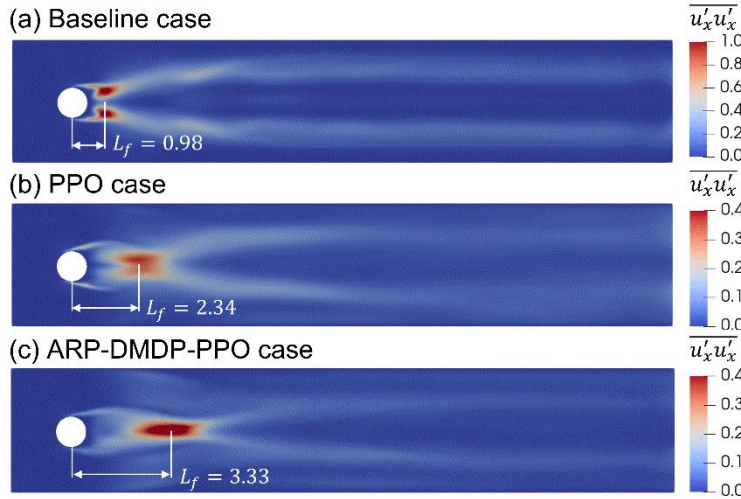
17 Figure 15 shows the streamwise  $\overline{u'_x u'_x}$  over the entire control process in the three cases. The location of peak  
 18 streamwise  $\overline{u'_x u'_x}$  gives of good indication of where the vortex formation zone ends. In the baseline case, the length  
 19 of vortex formation zone is  $L_f = 0.98$ . In the PPO controlled case, the vortex formation zone is extended to  $L_f =$   
 20  $2.34$  and the distance between the centres of vortex formation region and the central line of cylinder is  
 21 significantly reduced. Compared to the PPO case, the vortex formation length  $L_f$  in the ARP-DMDP-PPO case is  
 even greater and the vortex formation region of two side eddies almost overlap, resulting in a longer and steadier  
 recirculation zone.



1 In summary, in comparison to the PPO agent, the ARP-DMDP-PPO agent is more effective in extending the  
 2 length of the recirculation zone and moving the location of vortex shedding substantially downstream. This results  
 3 in a much steadier near wake flow and the significant reduction of  $\sigma_D$  and  $\sigma_L$  as presented in Table 2.



4  
 5 **Fig. 14.** Contours of time-averaged velocity over the entire control process in (a) Baseline case; (b) PPO case;  
 6 (c) ARP-DMDP-PPO case. Contours of time-averaged vorticity over the entire control process in (d)  
 7 baseline case; (e) PPO case; (f) ARP-DMDP-PPO case.



8  
 9 **Fig. 15.** Contours of  $\overline{u'_x u'_x}$  over the entire control process in (a) baseline case, (b) PPO case and (c) ARP-  
 10 DMDP-PPO case.

## 11 4. Conclusions

12 In this study, we have developed a new hybrid DRL method which combines a Markov decision process  
 13 (MDP) with time delays with the use of a first-order autoregressive policy (ARP), i.e. the ARP-DMDP-PPO



1 method. This hybrid DRL method is applied to control the vortex shedding process from a 2D circular cylinder  
2 using four synthetic jet actuators at a freestream Reynolds number of 400. The aim of this work is to improve the  
3 time series regression of ANN and control stability by reducing the large fluctuations in lift and drag  
4 accompanying the drag reduction; a problem which remains unsolved in the existing work and tends to escalate  
5 as the Reynolds number increases.

6 The Markov decision process (MDP) with time delays is introduced to quantify the action delays in the DRL  
7 process due to the time elapse between actuation and response of flow. A cross-correlation analysis has been  
8 performed between signals of the actuation jets and the velocity probes in the wake to validate the action delay in  
9 the decision process of DRL. The ARP-DMDP-PPO method has trained an agent successfully to implement a  
10 smooth and temporal coherent control. Compared to the standard DRL method, this method utilizes the historical  
11 samples without additional sampling in training and it is shown to be capable of reducing the magnitude of drag  
12 and lift fluctuations by approximately 90% while achieving a similar level of drag reduction in the deterministic  
13 control at the same actuation frequency.

14 In comparison to the existing the standard PPO method, this new method has yielded a stable and coherent  
15 control which results in a steadier and more elongated vortex formation zone behind the cylinder hence a much  
16 weaker vortex shedding process and less fluctuating lift and drag forces. This study demonstrates the necessity of  
17 including a physics-informed time delay in the MDP and the benefits of introducing ARPs to achieve a robust and  
18 temporal-coherent control of unsteady forces in active flow control. The pseudo-periodic control signals that the  
19 control strategy produced by the DRL algorithms are also more regular making them more applicable to AFC in  
20 real settings. Reduction of high-level temporal fluctuations in drag and lift will help to decrease dynamic loads  
21 and structural fatigue leading to an improved structural durability and operational safety which will benefit many  
22 industrial applications.

## 24 **Acknowledgement**

25 Yiqian Mao would like to acknowledge the Scholarship awarded from both the University of Manchester  
26 and the China Scholarship Council to his PhD studies. The authors would like to thank Dr. Jean Rabault  
27 (University of Oslo, Oslo, Norway) and Mr. Hongwei Tang (Nanjing University of Aeronautics and Astronautics,  
28 Nanjing, China) for making their open-source codes for deep reinforcement learning and numerical simulation  
29 available online at <https://github.com/jerabaul29/Cylinder2DFlowControlDRLParallel> and  
30 <https://github.com/thw1021/Cylinder2DFlowControlGeneral>.

## 1 **Author declarations**

## 2 **Conflict of interest**

3 The authors declare that they have no conflict of interest.

## 4 **Data availability**

5 The data that support the findings of this study are available from the corresponding author upon reasonable  
6 request.

## 7 **References**

- 8 <sup>1</sup> S. Zhong, M. Jabbal, H. Tang, L. Garcillan, F. Guo, N. Wood, and C. Warsop, "Towards the design of synthetic-  
9 jet actuators for full-scale flight conditions : Part 1: The fluid mechanics of synthetic-jet actuators," *Flow, Turbul.*  
10 *Combust.* **78**, 283-307 (2007).
- 11 <sup>2</sup> X. Li, F. Wu, Y. Tao, M. Yang, R. Newman, and D. Vainchtein, "Numerical study of the air flow through an  
12 air-conditioning unit on high-speed trains," *J. Wind Eng. Ind. Aerodyn.* **187**, 26 (2019).
- 13 <sup>3</sup> M. Urquhart, M. Varney, S. Sebben, and M. Passmore, "Drag reduction mechanisms on a generic square-back  
14 vehicle using an optimised yaw-insensitive base cavity," *Exp. Fluids* **62**(241), (2021).
- 15 <sup>4</sup> J.C. Schulmeister, J.M. Dahl, G.D. Weymouth, and M.S. Triantafyllou, "Flow control with rotating cylinders,"  
16 *J. Fluid Mech.* **825**, 743 (2017).
- 17 <sup>5</sup> H. Choi, W.P. Jeon, and J. Kim, "Control of flow over a bluff body," *Annu. Rev. Fluid Mech.* **40**, 113 (2008).
- 18 <sup>6</sup> C. Zhu, J. Chen, J. Wu, and T. Wang, "Dynamic stall control of the wind turbine airfoil via single-row and  
19 double-row passive vortex generators," *Energy* **189**, 116272 (2019).
- 20 <sup>7</sup> S.S. Collis, R.D. Joslin, A. Seifert, and V. Theofilis, "Issues in active flow control: theory, control, simulation,  
21 and experiment," *Prog. Aerosp. Sci.* **40**(4–5), 237 (2004).
- 22 <sup>8</sup> T. Shaqarin, P. Oswald, B.R. Noack, and R. Semaan, "Drag reduction of a D-shaped bluff-body using linear  
23 parameter varying control," *Phys. Fluids* **33**, (2021).
- 24 <sup>9</sup> B. Plumejeau, S. Delprat, L. Keirsbulck, M. Lippert, and W. Abassi, "Ultra-local model-based control of the  
25 square-back Ahmed body wake flow," *Phys. Fluids* **31**(8), 085103 (2019).
- 26 <sup>10</sup> D. Gao, H. Meng, Y. Huang, G. Chen, and W.L. Chen, "Active flow control of the dynamic wake behind a  
27 square cylinder using combined jets at the front and rear stagnation points," *Phys. Fluids* **33**(4), 047101 (2021).
- 28 <sup>11</sup> M. Gadalla, M. Cianferra, M. Tezzele, G. Stabile, A. Mola, and G. Rozza, "On the comparison of LES data-

1 driven reduced order approaches for hydroacoustic analysis," *Comput. Fluids* **216**, 104819 (2021).

2 <sup>12</sup> A. Towne, O.T. Schmidt, and T. Colonius, "Spectral proper orthogonal decomposition and its relationship to  
3 dynamic mode decomposition and resolvent analysis," *J. Fluid Mech.* **847**, 821 (2018).

4 <sup>13</sup> S.L. Brunton, B.R. Noack, and P. Koumoutsakos, "Machine Learning for Fluid Mechanics," *Annu. Rev. Fluid  
5 Mech.* **52**, 477 (2020).

6 <sup>14</sup> S. Shimomura, S. Sekimoto, A. Oyama, K. Fujii, and H. Nishida, "Closed-loop flow separation control using  
7 the deep q network over airfoil," *AIAA J.* **58**(10), 4260 (2020).

8 <sup>15</sup> Y. LI, J. CHANG, C. KONG, and W. BAO, "Recent progress of machine learning in flow modeling and active  
9 flow control," *Chinese J. Aeronaut.* **35**(4), 14 (2021).

10 <sup>16</sup> K. Fukami, K. Fukagata, and K. Taira, "Assessment of supervised machine learning methods for fluid flows,"  
11 *Theor. Comput. Fluid Dyn.* **34**(4), 497 (2020).

12 <sup>17</sup> P. Zwintzsch, F. Gómez, and H.M. Blackburn, "Data-driven control of the turbulent flow past a cylinder," *J.  
13 Fluids Struct.* **89**, 232 (2019).

14 <sup>18</sup> H. Xu, W. Zhang, and Y. Wang, "Explore missing flow dynamics by physics-informed deep learning: The  
15 parameterized governing systems," *Phys. Fluids* **33**(9), 095116 (2021).

16 <sup>19</sup> E. Kharazmi, D. Fan, Z. Wang, and M.S. Triantafyllou, "Inferring vortex induced vibrations of flexible cylinders  
17 using physics-informed neural networks," *J. Fluids Struct.* **107**, 103367 (2021).

18 <sup>20</sup> G. Minelli, T. Dong, B.R. Noack, and S. Krajnović, "Upstream actuation for bluff-body wake control driven by  
19 a genetically inspired optimization," *J. Fluid Mech.* **893**, A1 (2020).

20 <sup>21</sup> F. Ren, C. Wang, and H. Tang, "Active control of vortex-induced vibration of a circular cylinder using machine  
21 learning," *Phys. Fluids* **31**(9), 093601 (2019).

22 <sup>22</sup> G. Ortali, N. Demo, and G. Rozza, "A Gaussian Process Regression approach within a data-driven POD  
23 framework for engineering problems in fluid dynamics," *Math. Eng.* **4**(3), 1 (2022).

24 <sup>23</sup> N. Alhazmi, Y. Ghazi, M. Aldosari, R. Tezaur, and C. Farhat, "Training a neural-network-based surrogate  
25 model for aerodynamic optimization using a gaussian process", in *AIAA Scitech 2021 Forum* (AIAA, 2021), pp.  
26 1–9.

27 <sup>24</sup> I. Kang, K.H. Lee, J.H. Lee, and J.W. Moon, "Artificial neural network-based control of a variable refrigerant  
28 flow system in the cooling season," *Energies* **11**(7), 1643 (2018).

29 <sup>25</sup> W. Zhang, Y. Yu, Y. Qi, F. Shu, and Y. Wang, "Short-term traffic flow prediction based on spatio-temporal  
30 analysis and CNN deep learning," *Transp. A Transp. Sci.* **15**(2), 1688 (2019).

- 1 <sup>26</sup> H. Eivazi, H. Veisi, M.H. Naderi, and V. Esfahanian, "Deep neural networks for nonlinear model order  
2 reduction of unsteady flows," *Phys. Fluids* **32**(10), 105104 (2020).
- 3 <sup>27</sup> B.Z. Han and W.X. Huang, "Active control for drag reduction of turbulent channel flow based on convolutional  
4 neural networks," *Phys. Fluids* **32**(9), 095108 (2020).
- 5 <sup>28</sup> C. Lee, J. Kim, D. Babcock, and R. Goodman, "Application of neural networks to turbulence control for drag  
6 reduction," *Phys. Fluids* **9**(6), 1740 (1997).
- 7 <sup>29</sup> K. Fukami, T. Nakamura, and K. Fukagata, "Convolutional neural network based hierarchical autoencoder for  
8 nonlinear mode decomposition of fluid field data," *Phys. Fluids* **32**(9), 095110 (2020).
- 9 <sup>30</sup> M. Morimoto, K. Fukami, K. Zhang, A.G. Nair, and K. Fukagata, "Convolutional neural networks for fluid  
10 flow analysis: toward effective metamodeling and low dimensionalization," *Theor. Comput. Fluid Dyn.* **35**, 633  
11 (2021).
- 12 <sup>31</sup> J.Z. Peng, S. Chen, N. Aubry, Z.H. Chen, and W.T. Wu, "Time-variant prediction of flow over an airfoil using  
13 deep neural network," *Phys. Fluids* **32**(12), 123602 (2020).
- 14 <sup>32</sup> T. Nakamura, K. Fukami, K. Hasegawa, Y. Nabae, and K. Fukagata, "Convolutional neural network and long  
15 short-term memory based reduced order surrogate for minimal turbulent channel flow," *Phys. Fluids* **33**(2),  
16 025116 (2021).
- 17 <sup>33</sup> K. Hasegawa, K. Fukami, T. Murata, and K. Fukagata, in *Fluid Dyn. Res.* (2020), p. 065501.
- 18 <sup>34</sup> R.S. Sutton, *Reinforcement Learning: An Introduction*, Second edition. (The MIT Press, Cambridge, MA,  
19 2018).
- 20 <sup>35</sup> J. Rabault, F. Ren, W. Zhang, H. Tang, and H. Xu, "Deep reinforcement learning in fluid mechanics: A  
21 promising method for both active flow control and shape optimization," *J. Hydrodyn.* **32**(2), 234 (2020).
- 22 <sup>36</sup> F. Ren, H. bao Hu, and H. Tang, "Active flow control using machine learning: A brief review," *J. Hydrodyn.*  
23 **32**(2), 247 (2020).
- 24 <sup>37</sup> P. Garnier, J. Viquerat, J. Rabault, A. Larcher, A. Kuhnle, and E. Hachem, "A review on deep reinforcement  
25 learning for fluid mechanics," *Comput. Fluids* **225**, 104973 (2021).
- 26 <sup>38</sup> Y.-Z. Wang, Y.-F. Mei, N. Aubry, Z. Chen, P. Wu, and W.-T. Wu, "Deep reinforcement learning based synthetic  
27 jet control on disturbed flow over airfoil," *Phys. Fluids* **34**(3), 033606 (2022).
- 28 <sup>39</sup> J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal Policy Optimization Algorithms,"  
29 arXiv:1707.06347, (2017).
- 30 <sup>40</sup> S. Fujimoto, H. Van Hoof, and D. Meger, "Addressing Function Approximation Error in Actor-Critic Methods,"

1 in *35th Int. Conf. Mach. Learn. ICML 2018* (PMLR, 2018), pp. 1587–1596.

2 <sup>41</sup> D. Fan, L. Yang, Z. Wang, M.S. Triantafyllou, and G.E. Karniadakis, "Reinforcement learning for bluff body  
3 active flow control in experiments and simulations," *Proc. Natl. Acad. Sci. U. S. A.* **117**(42), 26091 (2020).

4 <sup>42</sup> F. Ren, C. Wang, and H. Tang, "Bluff body uses deep-reinforcement-learning trained active flow control to  
5 achieve hydrodynamic stealth," *Phys. Fluids* **33**(9), 093602 (2021).

6 <sup>43</sup> R. Paris, S. Beneddine, and J. Dandois, "Robust flow control and optimal sensor placement using deep  
7 reinforcement learning," *J. Fluid Mech.* **913**, A25 (2021).

8 <sup>44</sup> J. Li and M. Zhang, "Reinforcement-learning-based control of confined cylinder wakes with stability analyses,"  
9 *J. Fluid Mech.* **932**, A44 (2022).

10 <sup>45</sup> J. Dai, P. Liu, Q. Qu, L. Li, and T. Niu, "Aerodynamic optimization of high-lift devices using a 2D-to-3D  
11 optimization method based on deep reinforcement learning and transfer learning," *Aerosp. Sci. Technol.* **121**,  
12 107348 (2022).

13 <sup>46</sup> J. Rabault, M. Kuchta, A. Jensen, U. Réglade, and N. Cerardi, "Artificial neural networks trained through deep  
14 reinforcement learning discover control strategies for active flow control," *J. Fluid Mech.* **865**, 281 (2019).

15 <sup>47</sup> J. Rabault and A. Kuhnle, "Accelerating deep reinforcement learning strategies of flow control through a multi-  
16 environment approach," *Phys. Fluids* **31**(9), 094105 (2019).

17 <sup>48</sup> H. Tang, J. Rabault, A. Kuhnle, Y. Wang, and T. Wang, "Robust active flow control over a range of Reynolds  
18 numbers using an artificial neural network trained through deep reinforcement learning," *Phys. Fluids* **32**(5),  
19 053605 (2020).

20 <sup>49</sup> M.M. Zdravkovich, *Flow around Circular Cylinders: A Comprehensive Guide through Flow Phenomena,*  
21 *Experiments, Applications, Mathematical Models, and Computer Simulations* (Oxford University Press, Oxford,  
22 1997).

23 <sup>50</sup> F. Ren, J. Rabault, and H. Tang, "Applying deep reinforcement learning to active flow control in weakly  
24 turbulent conditions," *Phys. Fluids* **33**(3), 037121 (2021).

25 <sup>51</sup> Q. Zhang, C. Su, M. Tsubokura, Z. Hu, and Y. Wang, "Coupling analysis of transient aerodynamic and dynamic  
26 response of articulated heavy vehicles under crosswinds," *Phys. Fluids* **34**(1), 017106 (2022).

27 <sup>52</sup> X. Chen, T. Liu, Y. Xia, W. Li, Z. Guo, Z. Jiang, and M. Li, "The evolution of airtight performance for a high-  
28 speed train during its long-term service," *Measurement* **177**, 109319 (2021).

29 <sup>53</sup> S. Meng, D. Zhou, X. Xiong, and G. Chen, "The Effect of the Nose Length on the Aerodynamics of a High-  
30 Speed Train Passing Through a Noise Barrier," *Flow, Turbul. Combust.* 2021 1082 **108**, 411 (2021).

- 1 <sup>54</sup> K. Goda, "A multistep technique with implicit difference schemes for calculating two- or three-dimensional  
2 cavity flows," *J. Comput. Phys.* **30**(1), 76 (1979).
- 3 <sup>55</sup> A. Logg, K.A. Mardal, and G. Wells, *Automated Solution of Differential Equations by the Finite Element  
4 Method: The FEniCS Book (Lecture Notes in Computational Science and Engineering)* (Springer, 2012).
- 5 <sup>56</sup> A. Kuhnle, M. Schaarschmidt, and K. Fricke, Tensorforce: a TensorFlow library for applied reinforcement  
6 learning, <https://github.com/tensorforce/tensorforce>.
- 7 <sup>57</sup> K. V. Katsikopoulos and S.E. Engelbrecht, "Markov decision processes with delays and asynchronous cost  
8 collection," *IEEE Trans. Automat. Contr.* **48**(4), 568 (2003).
- 9 <sup>58</sup> M.S. Bloor, "The transition to turbulence in the wake of a circular cylinder," *J. Fluid Mech.* **19**(2), 290 (1964).
- 10 <sup>59</sup> G. Chopra and S. Mittal, "Drag coefficient and formation length at the onset of vortex shedding," *Phys. Fluids*  
11 **31**(1), 013601 (2019).
- 12 <sup>60</sup> M.S. Bloor and J.H. Gerrard, "Measurements on turbulent vortices in a cylinder wake," *Proc. R. Soc. London,*  
13 *Ser. A* **294**(1438), 319–342 (1966).
- 14 <sup>61</sup> D. Korenkevych, A. Rupam Mahmood, G. Vasan, and J. Bergstra, "Autoregressive policies for continuous  
15 control deep reinforcement learning," in *IJCAI Int. Jt. Conf. Artif. Intell.* (IJCAI, 2019), pp. 2754–2762.
- 16 <sup>62</sup> M.L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming* (John Wiley & Sons,  
17 2014).
- 18 <sup>63</sup> K. Liu, J. Deng, and M. Mei, "Experimental study on the confined flow over a circular cylinder with a splitter  
19 plate," *Flow Meas. Instrum.* **51**, 95 (2016).
- 20 <sup>64</sup> Y.Z. Law and R.K. Jaiman, "Wake stabilization mechanism of low-drag suppression devices for vortex-induced  
21 vibration," *J. Fluids Struct.* **70**, 428 (2017).

22