

A review on deep reinforcement learning for fluid mechanics: an update

J Viquerat, P Meliga, A Larcher, E Hachem

▶ To cite this version:

J Viquerat, P Meliga, A Larcher, E Hachem. A review on deep reinforcement learning for fluid mechanics: an update. 2021. hal-03432654v1

HAL Id: hal-03432654 https://hal.science/hal-03432654v1

Preprint submitted on 17 Nov 2021 (v1), last revised 17 Jan 2023 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A REVIEW ON DEEP REINFORCEMENT LEARNING FOR FLUID MECHANICS: AN UPDATE

A Preprint

J. Viquerat* MINES Paristech, CEMEF PSL - Research University jonathan.viquerat@mines-paristech.fr

> **A. Larcher** MINES Paristech, CEMEF PSL - Research University

P. Meliga MINES Paristech, CEMEF PSL - Research University

E. Hachem MINES Paristech, CEMEF PSL - Research University

October 20, 2021

Abstract

In the past couple of years, the interest of the fluid mechanics community for deep re-1 inforcement learning (DRL) techniques has increased at fast pace, leading to a growing 2 bibliography on the topic. Due to its ability to solve complex decision-making prob-3 lems, DRL has especially emerged as a valuable tool to perform flow control, but 4 recent publications also advertise great potential for other applications, such as shape 5 optimization or micro-fluidics. The present work proposes an exhaustive review of the existing literature, and is a follow-up to our previous review on the topic. The 7 contributions are regrouped by domain of application, and are compared together re-8 garding algorithmic and technical choices, such as state selection, reward design, time 9 granularity, and more. Based on these comparisons, general conclusions are drawn 10 regarding the current state-of-the-art, and perspectives for future improvements are 11 sketched. 12

13 Keywords Deep reinforcement learning · Fluid mechanics

14 **1** Introduction

During the past decade, machine learning methods, and more specifically deep neural network (DNN), have achieved great successes in a wide variety of domains. State-of-the-art neural network architectures have reached astonishing performance levels in image classification tasks [1, 2], speech recognition [3] or generative tasks [4]. With a generalized access to GPU computational resources through cheaper hardware or cloud computing, such advances have been paving the way for a general evolution of the reference methods in these domains at both academic and industrial levels.

The rapid expansion of neural networks to multiple domains has also yielded important progress in the domain of decision-making techniques, by the coupling of DNNs with reinforcement learning algorithms (called deep reinforcement learning, or DRL). Several major obstacles that had been hindering classical reinforcement learning have been lifted using the feature extraction capabilities of DNNs, and their

²⁵ ability to handle high-dimensional state spaces. Unprecedented efficiency has been achieved in many

^{*}Corresponding author



Figure 1: **Timeline of recent publications on deep reinforcement learning for fluid dynamics.** Colors indicate different fields of application. Please note that we retain here the date of the first pre-print publication, and not that of final publication in peer-review journals. Indeed, the fast-paced evolution of the DRL community brings particular importance to pre-prints, sometimes supported by code release. The square symbols denote references included in our previous review [17].

domains such as robotics [5], language processing [6], or games [7, 8], but DRL has also proven useful in many industrial applications, such as autonomous cars [9, 10], or data center cooling [11].

The efforts for applying DRL to fluid mechanics are ongoing but still at an early stage, with only a 28 handful of pioneering studies providing insight into the performance improvements to be delivered in 29 the field. Nonetheless, from a few limital contributions in 2016 [12] and early 2018 [13], the domain has 30 undergone an increasing inflow of contributions, that pinnacled in 2020, with no less than 16 pre-prints 31 and articles, and a clear focus on drag reduction problems, as shown in figure 1. This enthusiasm can 32 be explained by two main factors: first, the increasing number of open-source initiatives [14, 15, 16], 33 that has led to an accelerated diffusion of the methods in the community. Second, the sustained 34 commitment from the machine learning community, that has allowed concurrently expanding the 35 scope from computationally inexpensive, low-dimensional reductions of the underlying fluid dynamics 36 to complex Navier–Stokes systems, all the way to experimental set-ups. The present review proposes 37 a six-year perspective of deep reinforcement learning applied to fluid flow problems, in the context 38 of both numerical and experimental environments. It is intended as a follow up to our first review 39 released as a pre-print in 2019 [17], that was followed in 2020 by a shorter review by another group of 40 authors, focused on drag reduction and shape optimization problems [18]. To the best of the authors 41 knowledge, those are the only other similar initiatives preceding this one, that are also featured in 42 figure 1 for the sake of completeness. 43

The organization is as follows: a reminder on the main DRL algorithms that have been used in a fluid dynamics context is provided in section 2. Section 3 lists of relevant issues to consider when evaluating the progress of DRL for practically fluid flow problems. An extensive bibliographical review is then conducted in sections 4 and 5, that covers a total of 32 papers, most of them subsequent to the previous review articles. Contributions are grouped and compared by domain of applications. Finally, section 6 draws general conclusions on the state-of-the-art and proposes a transversal study including suggestions for future work in the field.

⁵¹ 2 Deep reinforcement learning

This section covers the basic concepts of deep reinforcement learning, and briefly describes the methods 52 most represented in the selected contributions. First, the basic concepts of reinforcement learning are 53 introduced, whereafter value-based and policy-based methods are distinguished. Then, specificities of 54 DRL are detailed, and a curated list of algorithms is proposed, including deep Q-networks (DQN), 55 advantage actor-critic (A2C), proximal policy optimization (PPO), trust-region policy optimization 56 57 (TRPO), deep deterministic policy gradients (DDPG), twin-delayed deep deterministic policy gradients (TD3) and policy-based optimization (PBO). For a more sophisticated introduction to DRL, the reader 58 is referred to [19]. The notations used in the remaining of this review can be found in table 1. 50

	-
γ	discount factor
$\dot{\lambda}$	learning rate
α, β	step-size
ϵ	probability of random action
$^{s,s'}$	states
\mathcal{S}^+	set of all states
${\mathcal S}$	set of non-termination states
a	action
\mathcal{A}	set of all actions
r	reward
${\cal R}$	set of all rewards
\mathcal{D}	set of collected transitions
t	time station
T	final time station
a_t	action at time t
s_t	state at time t
r_t	reward at time t
r(s,a)	reward received for taking action a in state s
R(au)	discounted cumulative reward following trajectory τ
G_t	discounted cumulative reward starting from time t
π	policy
θ, θ'	parameterization vector of a policy
$\pi_{ heta}$	policy parameterized by θ
$\pi(s)$	action probability distribution in state s following π
$\pi(a s)$	probability of taking action a in state s following π
$V^{\pi}(s)$	value of state s under policy π
$V^*(s)$	value of state s under the optimal policy
$Q^{\pi}(s,a)$	value of taking action a in state s under policy π
$Q^*(s,a)$	value of taking action a in state s under the optimal policy
$Q_{\theta}(s,a)$	estimated value of taking action a in state s with parameterization θ

Table 1: Notations used in the present review

60 2.1 Reinforcement learning

Reinforcement learning is a class of methods designed for decision-making problems, in which an agent learns to interact with an environment by (i) observing it, (ii) taking actions based on these observations, and (iii) receiving rewards from it, as a measure of the quality of the action taken. RL is based on Markov decision processes, for which a typical execution goes as follows (see also figure 2):

 \diamond Assume the environment is in state $s_t \in S$ at iteration t, where S is a set of states;

- ⁶⁶ \diamond The agent uses w_t , an observation of the current environment state (and possibly a partial ⁶⁷ subset of s_t) to take action $a_t \in \mathcal{A}$, where \mathcal{A} is a set of actions;
- \diamond The environment reacts to the action by transitionning from s_t to state $s_{t+1} \in \mathcal{S}$;
- ⁶⁹ \diamond The agent is fed with a reward $r_t \in \mathcal{R}$, where \mathcal{R} is a set of rewards, and a new observation ⁷⁰ w_{t+1} .

The steps described above repeat until a termination state is reached, and the succession of states and actions then define a finite trajectory $\tau = (s_0, a_0, s_1, a_1, ...)$. In any given state, the objective of the agent is to determine the adequate action to maximize its cumulative reward over an episode, *i.e.* over one trajectory. Most often, the quantity of interest is the discounted cumulative reward along a trajectory, defined as:

$$R(\tau) = \sum_{t=0}^{T} \gamma^t r_t \,, \tag{1}$$



Figure 2: RL agent and its interactions with the environment.

where T is the horizon of the trajectory (*i.e.* the terminal time station), and $\gamma \in [0,1]$ is a discount 76 factor that weights the relative importance of present and future rewards. Within the zoology of DRL 77 methods, we distinguish two categories, namely model-based and model-free algorithms. Model-based 78 method incorporates a model of the environment they interact with, and will not be considered in this 79 paper (the reader is referred to [19] and references therein for details about model-based methods). On 80 the contrary, model-free algorithms directly interact with their environment, and are currently the most 81 commonly used within the DRL community, mainly for their ease of application and implementation. 82 Model-free methods are further distinguished between value-based methods and policy-based methods 83 [19]. Although both approaches aim at maximizing their expected return, policy-based methods do 84 so by directly optimizing the parameterized policy, while value-based methods learn to estimate the 85 expected value of a state-action pair optimally, which in turn determines the best action to take in 86 each state. 87

88 2.1.1 Value-based methods

⁸⁹ In value-based methods, the agent learns to optimally estimate a *value function*, which in turn dictates

⁹⁰ the policy of the agent by selecting the action of the highest value. One usually defines the *state value* ⁹¹ *function*:

$$V^{\pi}(s) = \mathop{\mathbb{E}}_{\tau \sim \pi} \left[R(\tau) | s \right],$$

denoting the expected discounted cumulative reward starting in state s, then following trajectory τ according to policy π , and the state-action value function, or Q-function:

$$Q^{\pi}(s,a) = \mathop{\mathbb{E}}_{\tau \sim \pi} \left[R(\tau) | s, a \right],$$

94 denoting the same expected discounted cumulative reward starting in state s and taking action a. 95 Both values are quite obviously such that:

$$V^{\pi}(s) = \mathbb{E}_{a \sim \pi} \big[Q^{\pi}(s, a) \big],$$

⁹⁶ meaning that in practice, $V^{\pi}(s)$ is the weighted average of $Q^{\pi}(s, a)$ over all possible actions by the ⁹⁷ probability of each action. One of the main value-based methods in use is called Q-learning, as it relies

on the learning of the Q-function to find an optimal policy. In classical Q-learning, the Q-function

⁹⁹ is stored in a Q-table, which is a simple array representing the estimated value of the optimal Q-

function $Q^*(s, a)$ for each pair $(s, a) \in \mathcal{S} \times \mathcal{A}$. The Q-table is initialized randomly, and its values are

¹⁰¹ progressively updated as the agent explores the environment, until the Bellman optimality condition ¹⁰² [20] is reached:

$$Q^*(s,a) = r(s,a) + \gamma \max_{a'} Q^*(s',a'),$$
(2)

¹⁰³ at which point the Q-table estimate of the Q-value has converged, and taking the action with the ¹⁰⁴ highest Q-value systematically leads to the optimal policy.

105 2.1.2 Policy-based methods

Policy methods maximize the expected discounted cumulative reward of a policy $\pi(a|s)$ mapping states to actions, and resort not to a value function, but to a probability distribution over actions given states. Compared to value-based methods, policy-based methods offer three main advantages:

they have better convergence properties, although they tend to get trapped in local minima;
they naturally handle high dimensional action spaces;

 \diamond they can learn stochastic policies.

Most DRL algorithms applied to fluid mechanics problems are policy gradient methods, in which gradient ascent is used to optimize a parameterized policy $\pi_{\theta}(a|s)$ with respect to some measure of the expected return. In practice, one defines an objective function based on the expected discounted cumulative reward:

$$J(\boldsymbol{\theta}) = \mathop{\mathbb{E}}_{\boldsymbol{\tau} \sim \pi_{\boldsymbol{\theta}}} \big[R(\boldsymbol{\tau}) \big],$$

and seeks the optimal parameterization θ^* that maximizes $J(\theta)$:

$$\theta^* = \arg \max_{\theta} \mathop{\mathbb{E}}_{\tau \sim \pi_{\theta}} \left[R(\tau) \right],$$

which can be done on paper by plugging an estimator of the policy gradient $\nabla_{\theta} J(\theta)$ into a gradient ascent algorithm. In practice, this is no small task, as one is looking for the gradient with respect to the policy parameters θ , in a context where the effects of policy changes on the state distribution are unknown (since modifying the policy will most likely modify the set of visited states, which will in turn affect performance in some indefinite manner). The standard derivation relies on the log-probability trick [21], and allows expressing $\nabla_{\theta} J(\theta)$ as an evaluable expected value:

$$\nabla_{\theta} J(\theta) = \mathop{\mathbb{E}}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T} \nabla_{\theta} \log \left(\pi_{\theta}(a_t | s_t) \right) R(\tau) \right],$$
(3)

¹²³ after which the gradient is used to update the policy parameters:

$$\theta \leftarrow \theta + \lambda \nabla_{\theta} J(\theta). \tag{4}$$

124 2.2 Deep reinforcement learning

Deep reinforcement learning (or DRL) is the result of applying RL using deep neural networks to output either value functions (value-based RL methods), or action distributions given input states (policy-based RL methods)². A neural network (NN) is a collection of artificial neurons, *i.e.* connected computational units with universal approximation capabilities [22, 23], that can be trained to arbitrarily well approximate the mapping function between input and output spaces. Each connection

 $^{^{2}}$ An alternative presented above is to use tables to store the values for every state or state-action pair, but such a strategy generally does not scale with the size of state-action spaces, and is thus limited to discrete spaces.



Figure 3: Fully connected neural network with two hidden layers.

provides the output of a neuron as an input to another neuron. Each neuron performs a weighted sum 130 of its inputs, to assign significance to the inputs with regard to the task the algorithm is trying to 131 learn. It then adds a bias to better represent the part of the output that is actually independent of the 132 input. Finally, it feeds a non-linear activation function that determines whether and to what extent 133 the computed value should affect the ultimate outcome. As sketched in figure 3, a fully connected 134 network is generally organized into layers, with the neurons of one layer being connected solely to 135 those of the immediately preceding and following layers. The layer that receives the external data is 136 the input layer, the layer that produces the outcome is the output layer, and in between them are zero 137 or more hidden layers. 138

The learning process in neural networks consists in adjusting all the biases and weights of the network 139 in order to reduce the value of a well-chosen loss function that represents the quality of the network 140 prediction. This update is usually performed by a stochastic gradient method, in which the gradients 141 of the loss function with respect to the weights and biases (*i.e.* the parameterization θ of the neural 142 network) are obtained using a back-propagation algorithm. The abundant literature available on this 143 topic (see [24] and the references therein) points out that a relevant network architecture (e.g. type 144 of network, depth, width of each layer), finely tuned hyper parameters (*i.e.* parameters whose value 145 cannot be estimated from data, e.g., optimizer, learning rate, batch size) and a sufficiently large 146 amount of data to learn from are key ingredients for a successful network training. 147

148 2.3 DRL algorithms

This section briefly reviews some of the most popular DRL methods encountered in the field of DRL for fluid mechanics. Only their main features are reviewed here, the reader interested in further details (or in the numerous custom variations introduced in the RL literature) is referred to the various references given in this section.

153 2.3.1 Deep and double deep Q-networks (DQN and DDQN)

In Q-learning methods, obtaining a converged Q-table for large state and actions spaces can be particularly expensive in terms of environment interactions. To overcome this issue, the map $S^+ \times A \longrightarrow \mathbb{R}$ is represented by a neural network, called deep Q-network [25], tasked with providing an estimate of the Q-value for each possible action given an input state. To do so, the Q-network is trained on state-action-reward transitions obtained by interacting with the environment. The loss used for the training is classically obtained from the Bellman equation (2):

$$L(\theta) = \mathbb{E}_{(s,a,r,s')\sim\mathcal{D}}\left[\frac{1}{2}\left(\left[r(s,a) + \gamma \max_{a'} Q_{\theta}(s',a')\right] - Q_{\theta}(s,a)\right)^2\right],\tag{5}$$

where $Q_{\theta}(s, a)$ is the Q-value *estimate* provided by the DQN for action s and state a under network parameterization θ , and \mathcal{D} represents the set of transitions collected from the environment. The quantity $r(s, a) + \gamma \max_{a'} Q_{\theta}(s', a')$, denoted *target*, also appears in the Bellman equation (2), as the estimate $Q_{\theta}(s, a)$ is equal to the target (and $L(\theta)$ is thus zero) when the optimal set of parameters θ^* is reached.

In order to balance the trade-off between exploration and exploitation, the DQN algorithm classically 165 implements a stochastic exploration strategy called ϵ -greedy: before each action, a random parame-166 ter p is drawn in [0, 1] and compared to a user-defined value $\epsilon \in [0, 1]$, and the action prescribed by 167 $\max_a Q_{\theta}(s, a)$ is taken only if $p > \epsilon$ (otherwise a random action is taken). The value of ϵ usually de-168 creases during the learning process, thereby progressively reducing exploration in favour of exploitation. 169 Nevertheless, the performance of vanilla DQN remains limited. This has led to multiple developments 170 aimed at stabilizing learning and at improving performance, a handful of which have become standard 171 practice e.g., replay [26], target networks [25], or double Q-networks [27]. For the sake of clarity, we 172 shall not go into the specifics of these evolutions, for which the interested reader can instead refer to 173 [19] and the references therein. 174

175 2.3.2 Vanilla deep policy gradient

176 In policy methods, a stochastic gradient algorithm is used to perform network updates from the policy 177 loss:

$$L(\theta) = \mathop{\mathbb{E}}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T} \log \left(\pi_{\theta}(a_t | s_t) \right) R(\tau) \right].$$
(6)

whose gradient is equal to the policy gradient (3). The latter is computed with the back-propagation 178 algorithm with respect to each weight and bias by the chain rule, one layer at the time from the output 179 to the input layer. Such a method is also known as Monte Carlo policy gradient, as the loss (6) takes 180 the form of an expected value, that can be numerically calculated using an empirical average over a set 181 of full trajectories. However, if some low-quality actions are taken along the trajectory, their negative 182 impact will be averaged by the high-quality actions and will remain undetected. This problem can be 183 overcome using actor-critic methods, in which a Q-function evaluation is used in conjunction with a 184 policy optimization. 185

186 2.3.3 Advantage actor-critic (A2C)

¹⁸⁷ Different strategies are available to alleviate the high variance of training the agent from (6), for which ¹⁸⁸ it has become customary to replace the discounted cumulative reward by the *advantage function*:

$$A(s,a) = Q(s,a) - V(s),$$

that represents the improvement in the expected cumulative reward when taking action a in state s, compared to the average of all possible actions taken in state s. As a result, the loss function reads:

$$L(\theta) = \mathop{\mathbb{E}}_{\tau \sim \pi_{\theta}} \left[\sum_{t=0}^{T} \log \left(\pi_{\theta}(a_t | s_t) \right) A^{\pi_{\theta}}(s_t, a_t) \right].$$

In practice, the classical policy network (called *actor*) is used concurrently with a second network (called *critic*), that learns to predict the state-value function V(s). The advantage function is then approximated as

$$A(s_t, a_t) \sim r(s_t, a_t) + \gamma V(s_{t+1}) - V(s_t).$$

to avoid having a third network learn to predict the state-action value Q(s, a). In contrast to the Monte Carlo-style update of vanilla policy gradient methods, the actor-critic algorithm allows training the policy network in a temporal-difference manner, meaning that updates can be performed several times during an episode, thanks to the critic state-value estimate [28].

198 2.3.4 Trust-region and proximal policy optimization (TRPO and PPO)

The performance of policy gradient methods is hurt by the high sensitivity to the learning rate, *i.e.*, 199 the size of the step to be taken in the gradient direction. Indeed, small learning rates are detrimental 200 to learning, but large learning rates can lead to a performance collapse if the agent falls off the cliff 201 and restarts from a poorly performing state with a locally bad policy (an issue magnified by the fact 202 that the learning rate cannot be tuned locally). Trust region policy optimization (TRPO [29])) ensures 203 continuous improvement by leveraging second-order natural gradient optimization to update the policy 204 parameters within a trust-region of fixed maximum Kullback-Leibler divergence between previous and 205 current policies. Proximal policy optimization (PPO [30]) uses a simpler yet effective heuristic to 206 similarly avoid destructive updates. Namely, it relies on the clipped surrogate loss: 207

$$L(\theta) = \mathop{\mathbb{E}}_{(s,a) \sim \pi_{\theta_{\text{old}}}} \left[\min\left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_{old}}(a|s)}, g\left(\epsilon, A^{\pi_{\theta_{\text{old}}}}(s,a)\right) \right) A^{\pi_{\theta_{\text{old}}}}(s,a) \right],$$

208 where

$$g(\epsilon, A) = \begin{cases} (1+\epsilon)A & \text{if } A \ge 0, \\ (1-\epsilon)A & \text{if } A < 0, \end{cases}$$

and ϵ is the clipping range, a small, user-defined parameter defining how far away the new policy is allowed to go from the old one. The general picture is that a positive (resp. negative) advantage increases (resp. decreases) the probability of taking action *a* in state *s*, but always by a proportion smaller than ϵ , otherwise the min kicks in (2.3.4) and its argument hits a ceiling of $1 + \epsilon$ (resp. a floor of $1 - \epsilon$). This prevents stepping too far away from the current policy, and ensures that the new policy will behave similarly.

Due to its improved learning stability and its relatively robust behaviour with respect to hyperparameters, the PPO algorithm has received considerable attention in the DRL community. As shown in table 4, it is by far the most common DRL algorithm exploited in the context of DRL-based control for fluid dynamics.

219 2.3.5 Deep deterministic policy gradient (DDPG)

Deep deterministic policy gradient (DDPG) can be thought as a DQN algorithm for continuous actions spaces, that combines the learning of a Q-network $Q_{\theta}(s, a)$ (as in the DQN algorithm) and a deterministic policy network $\mu_{\phi}(s)$. As in DQN, the replay buffer and target network tricks are used, the latter being a key ingredient of the method. Looking back at the DQN loss (5), it is obvious that the max_{a'} $Q_{\theta}(s', a')$ term does not make sense in the context of a continuous action space. In DDPG, the latter is thus approximated using the target network, yielding the modified loss:

$$L(\theta) = \mathbb{E}_{(s,a,r,s')\sim\mathcal{D}} \left[\frac{1}{2} \left(\left[r(s,a) + \gamma \, Q_{\theta_{\text{targ}}}(s',\mu_{\phi_{\text{targ}}}(s')) \right] - \, Q_{\theta}(s,a) \right)^2 \right].$$
(7)

Hence, the policy $\mu_{\phi}(s)$ is expected to produce actions corresponding to a maximum value predicted by the Q-network, and therefore its loss is obtained straightforwardly as:

$$L(\phi) = \mathop{\mathbb{E}}_{s \sim \mathcal{D}} \left[Q_{\theta} \left(s, \mu_{\phi}(s) \right) \right].$$
(8)

Finally, a gaussian noise is usually applied to the predicted actions in order to achieve an efficient balance between exploration and exploitation.

230 2.3.6 Twin-delayed DDPG (TD3)

The Twin-delayed DDPG (TD3) algorithm is a refinement of the DDPG method that improves its learning stability and robustness against hyper-parameters [31]. For the sake of brevity, only the differences between the two methods are pointed out here, namely: the use of a second Q-network to avoid the common problem of overestimation of the Q-value, as is done in DDQN [27];

- 236 ◊ additional delays in the policy and target network updates;
- \diamond additional noises in the target actions.

Compared to standard DDPG, these three modifications largely improve the stability and performance
of the method. Yet, as shown in table 4, these two methods have received little attention in the field
of DRL-based control for fluid dynamics.

241 2.4 Single-step DRL

In several contributions assessed in this review, the optimal policy to be learnt by the neural network 242 is state-independent, as is notably the case in optimization and open-loop control problems. We group 243 here under the "single-step DRL" label the class of algorithms dedicated to this class of problems under 244 the premise that it may be enough for the neural network to get only one attempt per episode at finding 245 the optimal. In essence, the proposed methods inherit from deep policy gradient algorithms in the sense 246 that relevant probability density function parameters are obtained from neural networks trained using 247 a policy gradient-like loss. Yet, they also fall heir of evolutionary strategies (ES), as their successive 248 249 steps follow a generation/individual nomenclature, exploiting information from previous generations 250 in order to update the parameters of a probability density function. The seminal PPO-1 algorithm proceeds from the standard PPO algorithm (section 2.3.4) and samples actions isotropically from scalar 251 covariance matrices [32, 33, 16]. The follow-up policy-based optimization (PBO) algorithm relies on 252 a variant of the vanilla policy gradient method and delivers several major improvements by adopting 253 key heuristics from the covariance matrix adaptation evolution strategy (CMA-ES [34]), including the 254 use of a valid, full covariance matrix generated from neural network outputs [35]. 255

²⁵⁶ **3** Open challenges

Before delving into the specifics of the compiled papers, it is important to define a consistent list of 257 challenges to serve as a common thread to measure the progress of DRL in the context of fluid mechanics 258 applications (and also to examine the willingness of the community to take on these challenges). For 259 those challenges left mostly unanswered, section 6 proposes a series of possible mitigation strategies that 260 have received consideration in the literature, albeit in a different context. The retained challenges are 261 computational cost (more generally, sampling-efficiency), turbulence, robustness, partial observability, 262 delays, and any combination of them. Nonetheless, there are several other challenges that should be 263 considered on the second level to help bridge the gap between DRL capabilities and the requirements 264 of practical deployment, for instance multi-agent DRL (leveraging experience from multiple agents 265 learning concurrently) or multi-objective reward (training an agent in reasoning about several weighted 266 objectives), see, e.g. [36, 37, 38] for comprehensive domain-agnostic surveys. 267

Computational cost/sampling efficiency: the environment of computational fluid dynamics (CFD)
 problems is resource expensive, as it routinely involves numerical simulations with tens or hundreds of
 millions of degrees of freedom (unless an appropriate low-dimensional reduction is achieved, which in
 itself often proves very challenging). This is all the more problematic since classical RL methods have
 low sample efficiency, *i.e.* many trials are required for the agent to learn a purposive behavior.

Stochasticity (turbulence): most natural and engineered flows are turbulent and carry energy dis tributed over a wide range of scales with varying degrees of spatial and temporal coherence. Their
 dynamics therefore inherently includes some degree of stochasticity, which might lead to high variance
 gradient estimates that hamper learning.

Robustness: optimizing robust policies is a key issue for fluid flow applications, with multiple sources
of uncertainty relating to the occurrence of irregular transient dynamics and the high sensitivity to
initial conditions and system parameters variations (all non-normal amplification mechanisms associated with the asymmetry of the Navier–Stokes convection operator), not to mention the difficulty to
consistently ascertain the accuracy of the computed numerical solutions.

Partial observability: the traditional states space of fluid flow problems are easily prohibitively
 large for policy learning. The agent must therefore operate under partially observable environments,



Figure 4: **Different type of delays encountered in DRL environments.** *Pre-action delay* corresponds to the time delay existing between the moment of state collection and the moment when action is applied to the environment (this type of delay does not exist in the case of numerical environments). *Post-action delay* is defined as the time delay existing between the moment the action is applied to the environment, and the moment where it becomes fully effective in the dynamics of the system.

in which case the performance becomes highly dependent on the quality and relevance of the data
available for observation. This issue is strongly related to data-driven model reduction techniques for
large scale dynamical systems, which usually require using measures of observability as an information
quality metric.

Pre-action delays: in numerical environments, states are collected in the environment, provided to
 the agent, and actions are returned instantaneously, which amounts to artificially interrupting the
 lapse of time every time the agent must draw new actions. In real-world environments, a certain delay
 is inevitable due to data processing, data transition, and physical constraints of sensors and actuators,
 during which the environment keeps evolving, meaning that the agent actually takes actions based on
 out-dated states.

 \circ Post-action delays: an environment has an intrinsic response time that depends on the interplay 294 between transient amplification of the action-induced initial energy and non-linear saturation (the 295 former all the more important in fluid mechanics where non-normal systems are common occurrence). 296 This entails *post-action delays*, defined as the time interval between the moment an action is applied, 297 and the moment it efficiently reaches the current state, that can undermine the accuracy of the reward 298 estimation and even prevent learning if they exceed the Lyapunov time (the characteristic time scale on 299 which a dynamical system is chaotic). The two types of delays defined above are illustrated in figure 4, 300 the general picture being that post-action delays affect both numerical or experimental environment, 301 302 while pre-action delays affect only experimental environments (unless they are purposely included in a numerical model). 303

³⁰⁴ 4 DRL for computational fluid dynamics

Of the 32 papers compiled in the present review, 30 consider applying DRL to computational fluid dynamic (CFD) systems. Those are classified and presented here in one of the categories listed in table 2, to put similar papers in perspective with respect to one another and to point out their specificities.

308 4.1 Drag reduction

Drag reduction is by far the most represented application domain in the literature, with 12 different 309 papers implementing various control strategies using zero-mass-flow-rate jets [14, 40, 39, 43, 48, 46, 47], 310 rotating cylinders [41, 42, 45, 32], plasma actuators [44], or passive devices [32], as illustrated in 311 figure 5. Almost all studies focus on prototypal, two-dimensional (2-D) incompressible flows past 312 span-wise infinite cylinders (generally different sections of a single cylinder) subjected to a uniform 313 and/or parabolic velocity profile. As seen from the comparison between studies provided in table 3, 314 all DRL algorithms belong to the actor-critic category, with a clear preference for ready-to-use PPO 315 implementations (see section 2.3.4), either from Tensorforce [62], OpenAI baselines [63], or Stable 316 Baselines [64]. Regarding the CFD solvers, FeniCS [65] is well represented, mostly because the open-317 318 source diffusion of the seminal work from Rabault et al. [14] has been heavily re-used in follow-up works [40, 42, 43, 44, 45, 47]. Regarding the numerical implementation, since performing a relevant network 319

Category	Domain	Reference				
	Drag reduction	[39, 14, 40, 41, 42, 43, 44, 45, 46, 32, 47, 48]				
	Heat transfer	[49, 33]				
NT · 1	Microfluidics	[50, 51]				
Numerical	Swimming	[12, 13, 52, 53]				
	Shape optimization	[54, 55, 16, 56]				
	Other	[57, 15, 58, 59]				
	Drag reduction	[60]				
Experimental	Flow separation	[61]				
	Microfluidics	[51]				
Review	-	[17, 18]				

Table 2: Classification of the reviewed papers by domain of application. The most represented domain of application is drag reduction, with no less than 12 papers in total.







(c) Downstream rotating control cylinders.

(d) Symmetric plasma actuators.

Figure 5: Different drag reduction methods represented in the DRL literature, in the context of moderate Reynolds flows around a 2D circular cylinder. (5a) Zero-mass-flow-rate jets are used to blow or suck fluid on the lateral sides of the obstacle. There can be two or four, possibly tilted. (5b) An angular velocity is applied to the obstacle, in order to alter the downstream flow and reduce drag. (5c) Two small control cylinders, placed downstream of the obstacle, are given angular velocities in order to stabilize the shedding of the main cylinder. (5d) Two symmetric plasma actuators are controlled to alter the fluid flow near the flow-separation point, thus reducing the overall drag on the obstacle.

update requires evaluating a sufficient number of actions drawn from the current policy (which in turn
 requires computing the same amount of rewards from resource-expensive numerical simulations), most
 studies have the agent acquire experience at a faster pace by interacting with multiple environments

323 simultaneously. This has become a customary procedure after the methodological paper by Rabault

and Kuhnle [40] on the accelerated gathering of state-action-reward transitions, which highlighted an almost perfect speedup up to 20 parallel environments, and a decent performance improvement up to 60.

Regarding the flow regimes, almost all contributions assume laminar conditions with Revnolds numbers 327 in a range of one hundred to a few hundred. The only two exceptions are [48], where a weakly turbulent 328 case at an intermediate Reynolds number Re = 1000 is explicitly targeted, and [32], where moderately 329 large Reynolds number in the range of a few thousand to a few ten thousand are tackled in the frame 330 of Reynolds averaged Navier–Stokes (RANS); see figure 6 for an illustration pertaining to the fluidic 331 332 pinball, an equilateral triangle arrangement of rotating cylinders immersed in a turbulent stream. As stressed in [48], even weakly turbulent conditions make it significantly harder to achieve successful 333 drag reduction, as evidenced by the increased number of episodes needed to learn an efficient policy at 334

Control	Reference	Strategy	Re	DRL	CFD	n_{probes}	n_{act}	\mathbf{Actor}	$\delta t_{\mathbf{act}}/\delta t$	$\delta t_{\mathbf{act}}/\delta t_{\mathbf{phy}}$
	[14]	Jets	100	PPO (TFce)	FeniCS	151 (v)		[512, 512]	50	12
	[40]	*	100	PPO (TFce)	FeniCS	151 (v)		[512, 512]	50	12
	[39]	*	100	DDPG (IO)	UPACS	1(v)		[400, 300]	100	60
	[47]	*	100	PPO (TFce)	FeniCS	$63 \ (p)$		[512, 512]	50	12
	[46]	*	120	PPO-CMA (IO)	FastS	$3{-}16~(p)$		[512, 512]	50	22
Active/ Closed-loop	[43]	*	100 - 400	PPO (TFce)	FeniCS	236(?)	2	[512, 512]	200	30
	[48]	*	1000	PPO (IO)	LBM	$151 \ (v)$		ė	300	I
	[41]	Rotation	100	PPO (OpAI)	T-Flows	$12 \ (p)$		[64, 64]	30	20
	[45]	*	100 - 200	PPO (TFce)	FeniCS	476(p)	33	[512, 512]	85 - 350	10 - 20
	[42]	*	240	PPO (TFce)	FeniCS	99 (v)	5	ć	ć	ż
A attice /On and loom	[44]	Plasma	100	A2C (IO)	FeniCS	10(p)		[128, 64]	ż	
Acuve/Open-noop	[32]	Rotation	2200	PPO (StB1)	Cimlib	I	2^{-3}	[4, 4]	I	I
Passive	[32]	Device	100-22000	PPO (StBl)	Cimlib	I	2^{-3}	[4, 4]	I	I

luction DRL applications. Only explicitly stated informations were retained from the	on mark "?", while non-applicable data are noted with a double-dash "-". In the case where	with different parameter values, the most significant one was retained. $n_{\rm probes}$ corresponds	vation collection, while $n_{\rm act}$ represents the action dimensionality. The information in the	re used (in most contributions, the critic architecture is missing). Finally, the information	$t_{\rm act}/\delta t$ of the action to the simulation time-steps, and to the ratio $\delta t_{\rm act}/\delta t_{\rm phy}$ of the action	rtex shedding period).
nain features of drag reduction DRL applications. Only explicitly stat	ations are noted by a question mark "?", while non-applicable data are noted w	ed in the considered paper with different parameter values, the most significar	sed in the domain for observation collection, while $n_{\rm act}$ represents the action	nnected network architecture used (in most contributions, the critic architect	in respectively to the ratio $\delta t_{\rm act}/\delta t$ of the action to the simulation time-steps, i	e scale (here equal to the vortex shedding period).
Table 3: Summary of the	contributions. Missing inform	multiple studies were conduction	to the number of sensors pl_{δ}	actor column refers to fully-	in the last two columns perts	time-step to the physical tin



Figure 6: Vorticity field of the fluidic pinball case considered in [32]. The three cylinders are free to rotate at different angular velocities, leading to complex flow features in the near wake region.

higher Reynolds numbers. Moreover, the use of transfer learning from strategies learned at Re = 100335 to flow control at Re = 1000 is shown to be ineffective in this configuration, due to too different 336 flow dynamics. Nonetheless, it is shown possible in [43] to achieve robust flow control over a range 337 of Reynolds numbers by training simultaneously a single agent at four different Reynolds numbers 338 distributed between 100 and 400. After training, the agent succeeds in efficiently reducing drag for 339 Reynolds numbers in the range from 60 to 400, although the performance for each value of Re is 340 slightly lower than that achieved training an agent specifically at this Reynolds number. In [46], the 341 authors also underline the interest of using non-dimensionalized quantities as input states, which can 342 increase the robustness of control strategies even learned at a single Reynolds number. Yet, as stated 343 above about the work of [48], this approach is limited to cases with similar flow patterns, and does not 344 carry over to turbulent (not even weakly turbulent) flows. 345

About the numerical reward, most contributions rely on the design proposed in [14]:

$$r_t = -\langle C_D \rangle - \beta \left| \langle C_L \rangle \right|,\tag{9}$$

where the operator $\langle \cdot \rangle$ indicates the sliding average over one vortex shedding period T [14, 44, 46] or 347 over one action time-step δt_{act} [41, 48, 43, 45]. The parameter β varies in a range from 0.2 to 1 in 348 the papers reviewed in this section, and prevents the network to achieve efficient drag reduction by 349 relying on a large induced lift, as it is damageable in many practical applications. Whenever a different 350 reward is used, penalization terms associated with the cost of the control are rarely considered (with 351 reference [32] being an exception), as it has been found customary to explicitly bound the actuation 352 amplitude for the control to remain small compared to the system relevant physical quantities. Of 353 particular interest is the recent approach of [47] exploiting dynamic mode decomposition to design a 354 reward function based on mode amplitudes, that has led to efficient control strategies, although the 355 proposed approach supposes additional reward tuning compared to (9). 356

While the considered number of free action parameters $n_{\rm act}$ remains limited to 3 at most, the number 357 of probes used to collect observations varies considerably from one contribution to another, even for 358 similar setups. The baseline configuration consists of a certain number of velocity or pressure probes, 359 uniformly distributed in the vicinity of the cylinder as well as in its wake region, as illustrated in figure 360 7. The sensitivity of the learned control strategy to the probes distribution has been briefly studied in 361 [14]. A subsequent, more complete analysis has been performed in [46], where the authors evidence a 362 critical impact on the control performance, the information provided by sensors positioned in the near 363 wake being reportedly more relevant for learning than that of sensors positioned further downstream. 364 This can be attributed to the fact that, by observing the flow closely downstream of the cylinder, the 365 agent is able to observe the consequences of its actions right after they were taken, while more distant 366 sensors provide a delayed feedback that can be more difficult to interpret. To circumvent this issue, 367 the authors in [46] introduce a specific method, called sparse PPO-CMA, in which an optimal set of 368 sensors is automatically selected during the learning process. Another notable approach is that of [39], 369 370 where only one probe is used downstream of the cylinder, collecting observations at a higher rate than 371 the action time-step δt_{act} , and stacking them into a single observation vector when feeding them to the agent. In all relevant contributions, pressure or velocity are used indifferently as observations. 372



Figure 7: **Typical probe array for observation collection** in the context of drag reduction application on a 2D cylinder at moderate Reynolds numbers. Although the amount and position of probes vary, many contributions position probes in the vicinity of the obstacle and downstream of it. Insights on the impact of this choice regarding the performance of the agent can be found in [14] and [46].

As stated previously, the frequency at which the agent is allowed to provide new actions to the en-373 vironment is defined by an action time-step δt_{act} , that must be larger than the numerical simulation 374 time-step δt (for the agent to be able to observe the effects of its actions on the environment), but 375 smaller than the characteristic time scale $\delta t_{\rm phy}$ of the physical process to be controlled (for the actions 376 taken to be able to significantly alter the flow dynamics). Considerations about numerical stability 377 and accuracy of the numerical flow solution call for $\delta t \ll \delta t_{\rm phy}$, meaning that the user has considerable 378 leeway to adjust the action time step within these two bounds. In the reviewed contributions, the ratio 379 of the action time-step to the physical time scale $(\delta t_{\rm act}/\delta t_{\rm phy})$ is in a range of a few tens (*i.e.*, a few 380 tens of actions are taken per vortex shedding period). Meanwhile, the ratio of the action time-step to 381 the simulation time-step $(\delta t_{act}/\delta t)$ varies considerably with the Reynolds number, as it is set to a few 382 tens at Re = 100 [14, 41, 46], but increases up to a few hundreds at higher Reynolds values [43, 48, 45]. 383 Between each interaction with the environment, an interpolation scheme is usually exploited to avoid 384 abrupt control changes in the environment, which could cause numerical instabilities. A simple lin-385 ear interpolation is generally used [43], as the exponential decay initially considered in [14] has been 386 subsequently found to yield significant discontinuities in the lift computation. 387

Finally, the agent architecture networks are also consistent between the different studies, with two fully-388 connected layers used in all cases. While seven contributions appear to use pretty large layer sizes [14, 389 39, 43, 45, 46, 47], smaller networks are successfully used for problems of similar dimensionalities [41, 390 44]. To the best knowledge of the authors, no large-scale study of the impact of the network architecture 391 on the final agent performance was performed, and this choice remains most often empirical. Of 392 particular interest is the use of a tailored single-step method in [32], that allows performing passive 393 control with extremely small networks (see section 2.4), which is because the agent is not required to 394 learn a complex state-action relation, but only a transformation from a constant input state to a given 395 action. 396

397 4.2 Conjugate heat transfer

Although conjugate heat transfer systems governed by the coupled Navier–Stokes and heat equations seem natural candidates to extend the scope the DRL methodology and to increase the complexity of the targeted applications, the field has received little initial attention from the community. However, it may be starting gaining ground with two studies of DRL-based thermal control over the past two years.

In [49], the authors consider the closed-loop control of natural convection in a 2-D Rayleigh-Bénard convection cell simulated with an in-house lattice-Boltzmann code at Rayleigh numbers (based on the time-averaged temperature difference between the upper and lower plates) ranging from $Ra = 10^3$ (just before the onset of convection) to 10^7 (mild turbulence). The set-up, synthesized in figure 8, is as follows: the upper plate and the time-averaged lower plate temperature distributions are assumed constant. A discrete PPO agent whose implementation relies on OpenAI's stable-baselines [64] is



Figure 8: Illustration of the Rayleigh-Bénard convection control setup as presented in [49]. The top wall temperature is set to a constant temperature $T = T_c$, while the bottom temperature profile is cut in 10 segments on which the temperature can take values equal to $T_h + C$ or $T_h - C$. On the left and right walls, adiabatic conditions are imposed. Finally, the temperature and velocity fields are collected on a grid of probes equispaced in the computational domain.

then used to provide (after a normalization step) a zero-mean, piecewise-constant lower temperature fluctuation with the intent to reduce the convective effects. The actor is a fully-connected neural network with two hidden layers of width 64, and the instantaneous reward is defined as the opposite of the instantaneous Nusselt number Nu, which spurs the agent to minimize the convective effects at play. As in drag reduction applications, an array of probes is uniformy distributed over the computational domain to collect observations, under the form of both the temperature and velocity fields.

A particularity of this implementation is the systematic use of the four most recent observations in 415 the state buffer passed to the agent, an approach similar to that of [39], although the authors do 416 not provide insights about the impact of this choice on the agent performance. The agent is able 417 to entirely stabilize the convective flow up to $Ra = 10^5$, and consistently outperform state-of-the-art 418 linear approaches (proportional and proportional-derivative controllers) up to $Ra = 10^7$. Finally, the 419 authors also illustrate the controllability limits of the system using the simplified Lorenz attractor 420 system. By introducing a tunable artificial delay in the control, they show that exceeding half the 421 Lyapunov time in delay results in a highly degraded performance of the learned control. 422

Passive control of a similar Rayleigh-Bénard natural convection problem is performed in [33] with 423 the single-step approach presented in section 2.4. Compared to [49], the authors report excellent 424 control efficiency using much smaller networks (two hidden layers of width 2 vs. 64) and less parallel 425 environments (8 vs. 512) at $R_a = 10^4$, a value for which the optimal control determined in [49] ends up 426 being actually time-independent (unlike at higher Rayleigh numbers). The authors then use the same 427 approach and network architecture to minimize open-loop the inhomogeneity of temperature gradients 428 429 across the surface of two and three-dimensional hot workpieces under impingement cooling in a closed 430 cavity, identifying either optimal positions for cold air injectors relative to a fixed workpiece position, or optimal workpiece position relative to a fixed injector distribution (see illustration in figure 9). 431



Figure 9: Passive control of 3D forced convection in the context of workpiece cooling, reproduced from [33]. The positions of the three injectors are optimised in order to minimize the local temperature gradients in the workpiece during the cooling process.

432 4.3 Shape optimization

Shape optimization is another field fundamentally interrelated with flow control, that can seem as a 433 natural domain application for the DRL techniques covered above. Nonetheless, it is worth noticing 434 that shape optimization generally consists in determining a fixed shape meeting a set of required criteria 435 (e.q. high lift-to-drag ratio, low pressure loss). This is not per se the original purpose of DRL, that aims 436 at identifying optimal state-to-action relations (by means of neural network training) and is thus best 437 suited to dynamically manipulate a deformable shape. Two approaches exist in the literature in the 438 context of DRL-based shape optimization, a first one that directly optimizes state-independent shape 439 parameters (hence, direct shape optimization [16]) and a second one that incrementally modifies an 440 initial shape into an optimal one (hence, *incremental shape optimization* [54, 55, 56]). The conceptual 441 differences between these two approaches are illustrated in figure 10, and their implementations are 442 detailed in the following paragraphs.³ 443



Figure 10: Direct and incremental DRL-based shape optimization techniques present in the literature. In direct shape optimization (left), the agent is used as a proxy to optimize a direct mapping from a constant, initial state vector s_0 to the optimal state s^* , using a degenerate, single-step DRL algorithm. In incremental shape optimization (right), the agent learns the adequate mapping from the current state vector s_i to an incremental modification to apply to the latter, hence determining a path of incremental deformations to apply from s_0 to s^* . In both cases, multiple episodes are required for the agent to converge.

In direct shape optimization, the agent is used as a proxy to optimize a direct mapping from a constant, initial state vector s_0 to the optimal state s^* . This approach is implemented in [16] using single-step DRL (the degenerate class of DRL algorithms intended to optimize state-independent agent behavior) to design 2-D aerodynamic profiles without any *a priori* knowledge, feeding systematically an initial circle as input to the agent in single-step episodes (hence the adjective *stateless*). In practice, the shapes are described by a set of Bézier curves connecting the same number of control points, each with

³Although not directly included in the scope of the current review, it is worth mentioning the work from Lampton *et al.* [66], who considered the use of standard Q-learning method for shape optimization in 2008. In this contribution, optimization of airfoil geometries with four free parameters is considered, and the optimal policy is obtained by updating a Q-table in a temporal-difference fashion.



Figure 11: Shapes of optimal lift-to-drag ratio obtained with direct shape optimization with 3 free parameters (left), 9 free parameters (center) and 12 free parameters (right), reproduced from [16]. The shape parameterization relies on Bézier curves joining control points, the agent controlling their position and local curvature radius. These shapes were obtained by learning an optimal mapping from a simple cylinder.

⁴⁵⁰ 3 free parameters (2 coordinates, plus a local curvature radius). As shown in figure 11, the agent is ⁴⁵¹ able to design airfoil-like shapes maximizing the lift-to-drag ratio at Reynolds numbers of about a few ⁴⁵² hundred, which takes between one and three thousand CFD evaluations (*i.e.* single-step episodes) for ⁴⁵³ problem dimensionality ranging from 3 to 12, respectively.

The literature proposes three other DRL-based shape optimization contributions conversely relying on 454 incremental shape transformations, with the incremental modifications in [54, 56] taking the current 455 geometric parameters as input (of dimension 8 and 10, respectively), while the input states in [55] 456 consist in a distribution of wall Mach number (of dimension 4). In all three contributions, a pre-457 trained surrogate or a simplified model is used, either for full agent training, or to perform an initial 458 learning phase before re-training on a CFD environment using transfer learning. A key difference lies 459 in the fact that the authors in [54, 56] always use the same input state and consequently produce a 460 single optimized shape per training, while [55] relies on a set of input states randomly selected at the 461 beginning of each episode, meaning that the trained agent can be successfully re-used in production 462 on out-of-training input shapes. 463

From an algorithmic point of view, the choices are in line with those reported in the previous sections. 464 All algorithms are actor-critic, either PPO [55] or DDPG [54, 56], using fully-connected networks with 465 2 or 3 layers of width from 200 to 512 neurons per layer, except for [16]. As it represents an arbitrary 466 design choice in this specific application, the number of steps per episode is low, ranging from 5 to 467 20. A simple reward signal based on the lift-to-drag ratio is used in [54] and [16], but more complex 468 designs were used in the other two contributions. In [55], the instantaneous reward is based on the 469 difference of drag between the current and the previous generation, while [56] exploits a complex reward 470 expression based on the results of a principal component analysis. Overall, it is extremely difficult to 471 draw conclusions from these different approaches. Moving forward, a careful performance comparison 472 between direct and incremental approaches constitutes a topic of outmost importance, but a more 473 specific study focused on reward design could also be of great practical interest. 474

475 4.4 Swimming

The control of swimmers has been a pioneering field for applying deep reinforcement learning to fluid 476 mechanics problems, with a couple of contributions [12, 13] building on early seminal studies focusing 477 RL for schooling [67, 68]. In [12], the kinematics of two swimmers in a leader-follower configuration are 478 analyzed based on 2-D simulations of viscous incompressible flows. The first fish (leader) swims with 479 a steady gait and the second fish (follower) uses DRL to adapt its behaviour dynamically to account 480 for the effects of the wake encountered. The retained algorithm is DQN (see section 2.3.1), with input 481 482 states made up of the lateral displacements and orientation of the follower compared to leader, as 483 well as the two most recent actions, and the tail-beat status. An ϵ -greedy strategy is used to perform exploration, with randomness decaying from 0.5 to 0.1 over the course of learning. The reward design 484



Figure 12: Coordinated schooling of three swimmers, reproduced from [13]. The two followers interact with both rows of the wake shedding to increase their swimming efficiency.

is straightforward, and increasingly penalizes the follower when it strays too far away from the leader
 path:

$$r_t = 1 - 2\frac{|\Delta y|}{L},\tag{10}$$

where Δy is the aforementioned deviation, and L is the length of the swimmer. It takes roughly 100 000 transitions to learn the optimal behavior, and the results indicate that swimming in synchronized tandem (with the follower seeking to maintain its position in the center of the leader's wake, and its head synchronized with the vortices shed by the leader) can yield up to about 30% reduction in energy expenditure for the follower.

Reference [13] is a follow-up of [12] extended to 3-D schooling configurations, as illustrated in figure 492 12. A key contribution of this study is the use of a recurrent neural network, as the authors advertise 493 (and demonstrate by providing performance comparisons with standard feedforward neural network) 494 a greatly accelerated learning process using long-short term memory (LSTM) cells to encode the 495 unsteadiness of the value function, which in turn is found to enable far more robust smart-swimmers. 496 The retained recurrent network is composed of three layers of fully-connected LSTM units. The DQN 497 algorithm with Adam optimizer is used to perform training in a temporal-difference manner, using 498 again an ϵ -greedy exploration, with randomness decaying from 1 to 0.1. The training procedure requires 499 46 000 transitions (a reduction by roughly 50% with respect to the LSTM-less 2-D case). The results 500 support the conjecture that swimming in formation is energetically advantageous, with the trained 501 fishes showing collective energy-savings behaviors by appropriately placing themselves in appropriate 502 locations in the wake of other swimmers and interacting judiciously with their shed vortices. An almost 503 504 identical set-up (*i.e.* DQN algorithm exploiting an LSTM-based agent) is used in [52], to tackle a series 505 of different swimming problems, namely (i) point-to-point travel in quiescent flow, with reward based on normalized distance to target, (ii) holding a steady position in a rotating fluid flow, with reward 506 based on averaged translation velocity of the fish center of mass, and (iii) holding a steady position 507 in a Karman vortex street. The authors also emphasize the necessity to provide richer information 508 to the agent to reduce variability over multiple episodes. The retained approach consists in feeding 509 the agent with informations about the fish dynamics over the last four periods (e.g., depending on 510 the case, distance to objective, orientation of the swimmer, mean swimming velocities) and to add the 511 512 actions taken over the same history steps, which indeed is found to yield stable learning and efficient swimming strategies. 513

514 4.5 Microfluidics

⁵¹⁵ Micro-fluidics is one of the first fluid dynamics problems tackled with deep reinforcement learning ⁵¹⁶ techniques, but the related literature has since stalled to a single contribution from 2019 [50] (along ⁵¹⁷ with an additional experimental study [51] reviewed in section 5). In [50], the authors consider the



Figure 13: Rigid body control using steerable fluid jets, reproduced from [57]. Here, the goal is to push the ball back and forth from one jet to another.

inverse design problem of flow sculpting, in which a relevant sequence of micro-pillars is designed 518 to controllably deform an initial flow field into a desired one. A double-DQN agent (section 2.3.1) 519 is used that implements a convolutional policy, the full flow map being passed as input state [27]. 520 The agent network is composed of three convolutional/max-pooling layers followed by three batch-521 norm/fully-connected layers. The DDQN is supplemented with an experience replay method [69]. The 522 implemented reward is based on a pixel match rate (PMR) that measures the similarity of the current 523 flow with the target flow. This contribution also contains an interesting analysis comparing DDQN 524 performance with that of canonical methods, e.g., genetic algorithms and brute force approaches. 525

526 4.6 Other applications

This section connects to other contributions of the literature applying DRL to more restricted domains, *e.g.*, turbulence model generation [59], sloshing suppression [58] or instability mitigation in fluids [15], among others. It is worth insisting that the scarcity of publications on these topics does not reflect a lack of interest or priority, but rather the suddenness with which DRL has opened up new opportunities for a wide range of applications, as was already clear from the previous sections.

532 4.6.1 Flow control

In an early contribution by Ma et al. in 2018 [57], a TRPO (section 2.3.4) agent learns to play 533 different games (from rigid body balancing to complex music-playing games) based on the control of 534 rigid body by steerable fluid jets, as illustrated in figure 13. Regarding the environment, the Navier-535 Stokes equations are marched in time using a grid-based fluid-solid solver with adaptive refinement. A 536 convolutional auto-encoder trained on-the-fly is used to efficiently extract fluid flow features from the 537 environment. After their dimensionality has been reduced to an acceptable range, those are combined 538 with rigid body features and serve as input for the agent, which is shown to significantly improve the 539 learning speed compared to using rigid body features only. The state vector, whose size is lower than 540 100 elements, is fed to a standard fully-connected network of size [128, 64, 64, 32], which yields typical 541 training times in a range from 2 to 20 hours, depending of the game played. 542

Another under-represented type of application is the control of sloshing in tanks, despite obvious 543 practical interest for engineering applications, such as liquid carriers in ground, marine, or air transport 544 vehicles, as well as in earthquake excited water supply towers. Reference [58] is the only contribution 545 in the field, that considers suppressing sloshing in a tank initially submitted to a sinusoidal excitation 546 using two active controlled horizontal baffles. The comparison of two policy-gradient algorithms, 547 namely PPO (section 2.3.4) and TD3 (section 2.3.6) is a key contribution of this study. The state 548 information consists of the positions of the baffles, as well as the elevation and vertical velocities of two 549 additional probes in the tank. Given such inputs, the agent provides in return the horizontal velocities 550 to be applied to the baffles. For both algorithms, the actor is composed of a fully-connected network 551 with two layers of width 64. Actions are taken by the agent every 30 numerical time-steps, one episode 552 consisting in 200 actions, linearly interpolated from one time-step to the following. The reward is 553 equal to the time-averaged sloshing height, plus a penalization term to limit the displacements of the 554 baffles. Good convergence is reported both for PPO and TD3, although learning proves to be more 555 stable using TD3, as shown in figure 14. With direct learning, the authors notice a lack of robustness 556 when applying the learned strategy beyond the largest time used during training, which they show can 557 be overcome using behavior cloning to pre-train the agent. 558

Another noteworthy contribution is that of Belus *et al.* [15], that introduces a technique based on invariants intended for problems with large dimensional (up to 20) actions spaces. In this study, a



Figure 14: Performance comparison of PPO and TD3 for sloshing suppression task, reproduced from [58]. Although similar performance levels are obtained, the learning process proves to be more stable for TD3.

PPO agent (section 2.3.4) is used to mitigate the natural instabilities developping in a 1-D falling 561 liquid film using small jets blowing orthogonally to the flow direction. The number of jets and their 562 positions can vary, leading to different levels in control complexity. A three-layer, fully-connected 563 network of size [128, 64, 64] is used, with actions provided every 50 numerical time-steps to a variable 564 number of jets, based on local inputs recorded in the vicinity of each jet. Finally, the reward function 565 steers the agent to alleviate the waves arising from the instability of the flow. Three training methods 566 are compared, that differ by their ability to handle a large number of control jets: (i) local states are 567 concatenated and flattened before being fed to the actor, its output dimensionality being equal to the 568 number of jets; (ii) a similar approach is used, but instead of being flattened, the input states are fed 569 as is to a convolutional network; (iii) the vicinity of each jet is considered a local environment and used 570 to provide some states and a reward to a unique agent. This latter approach relies on the translational 571 invariance of the physical problem. It significantly enhances the amount of experience collected by 572 the agent during an episode, which the authors show allows tackling large dimensional action spaces 573 without increasing the amount of simulation time, as shown in figure 15. 574

575 4.6.2 Turbulence modeling

An approach somehow similar to that in [15] is used in another study by Novati *et al.* [59] to adjust 576 the coefficients of an eddy viscosity closure model in the attempt to reproduce the energy spectrum of 577 DNS computations. To this end, multiple agents are dispatched in the computational domain, with 578 each agent controlling locally the dissipation coefficient of the Smagorinsky SGS model. The provided 579 states are a blend of local (invariants of the gradient and Hessian of the velocity field) and global 580 quantities (modes of the energy spectrum, rate of viscous dissipation, total dissipation). Two types 581 of reward are proposed, based either on the Germano identity, or on a distance to a pre-computed 582 DNS spectrum. The agent uses the remember-and-forget experience replay method, in which networks 583 update are performed within a trust region, using a buffer holding the most recent transitions collected 584 by the policy (which supposedly greatly improves the sample efficiency by enabling data to be reused 585 586 multiple times for training) while dismissing those actions too unlikely under the current policy. The 587 network parameters are shared between agents, and their aggregated experiences are collected in a shared dataset used for training. 588



Figure 15: Learning curves obtained with the naïve learning technique (left) and with the invariance-based approach (right), reproduced from [15]. On the left figure, it can be seen that increasing the number of jets significantly increases learning time, here counted in number of actions. On the right figure, similar training times (counted in simulation steps) are required, whatever the number of jets (and therefore the action space dimension) used to control the instability.

589 5 Experimental fluid dynamics

The coupling of DRL and experimental fluid mechanics remains insufficiently explored, with only 3 out of the 32 papers compiled in this review applying DRL for experimental flow control purposes. Besides the possibly limited access to experimental devices for DRL practitionners, this is likely because several challenges such as controllability (the ability to efficiently reach a given state), observability (the ability to reliably measure changes in the state), sensitivity (to noise and system uncertainty) and system delays (see section 3) become increasingly important in experimental setups, even though they have received little attention in the context of idealized numerical environments.

597 5.1 Drag reduction

In [60], a drag reduction problem similar to that in figure 5c is considered, where an agent is given 598 control over the angular velocity of two rotating cylinders located in the wake of a fixed principal 599 cylinder. The Reynolds number is about $Re = 10^4$, and the agent is allowed to interact with the 600 601 environment every 0.1 s. An entire episode last 40 s, plus additional time consumption for initialisation 602 (4s, the time needed to wash out the transient before collecting any data) and for the reset procedure (2 mn, the time needed for the entire system to come back to rest). Overall, an experimental episode 603 lasts between 3 and 4 mn. A TD3 agent (section 2.3.6) based on Tensorflow is used, the updates being 604 performed only between episodes with a reward function similar to (9). The states provided to the 605 agent are the drag and lift coefficients measured on the main cylinder and the two control cylinders (an 606 approach noticeably different from that described in section 4.1). A key outcome of this study is the 607 necessity to high pass filter the experimental states before they are fed to the agent, as a comparison of 608 the performance with and without providing beforehand the experimental states as input to a Kalman 609 filter shows that the agent is essentially unable to learn an efficient strategy without the filtering stage. 610 Additional experiments are also performed to account for the power loss due to the friction of the 611 control cylinders 612

613 5.2 Flow separation

In [61], a DQN agent (section 2.3.1) learns to perform flow reattachment behind a NACA0015 airfoil by controlling the burst frequency of a plasma actuator at $Re = 6.3 \times 10^4$. Two different angles of attack are considered, namely 12° and 15°. The states provided to the agent consist of the unfiltered time-series data of the pressure at the surface of the airfoil, recorded through a set of 29 holes with high-frequency sensors, eventually downsampled to a total of 80 values. The actions are selected among a set of pre-defined burst frequencies, that includes four different values as well as an "off" choice. The reward is zero if the flow is not attached, and one if it is attached, as determined from the pressure coefficient at the trailing edge of the airfoil. The DQN agent achieves a satisfactory learning at the first angle of attack of 12°, with efficient strategies available after as little as 200 episodes, although not more efficient that a naive open-loop control with adequately selected burst frequency. Conversely, the agent significantly outperforms the naive open-loop design at the second angle of 15°, but learning is then much more challenging and takes about up to 800 episodes.

626 5.3 Microfluidics

The problem considered in [51] relates to the performance of microfluidics experiment platforms when 627 operated on extended periods of time. To overcome degraded flow stability beyond a certain timescale. 628 the authors introduce a DRL agent to adjust the flow conditions and maintain the experiment operabil-629 ity in an experimental device. Two low-Reynolds applications are considered, namely the positioning 630 of an interface between two miscible flows, and the dynamic control of the size of water-in-oil droplets 631 within a segmented flow. On both applications, the performances of a DQN [25] agent and a model-632 free episodic control (MFEC) [70] are compared, although it must be noted that the algorithm run 633 with different interaction frequencies (250 actions per episode for DQN, vs. 150 for MFEC) due to 634 equipment limitations. Observations are obtained from a high-speed camera and processed into an 635 84×84 pixels frame. In the first experiment, the reward is obtained calculating the distance between 636 the current observed interface and its target position, while in the second experiment it is computed 637 from the estimated radii of the generated droplets. The authors find that DQN requires a considerable 638 amount of frames (approximately 145000 in the first experiment) to surpass human-level performance, 639 albeit with large-scale fluctuations, while MFEC equires a reasonable number of frames to improve 640 641 and reach a stable level of performance (approximately 11 000 frames in the first experiment), but does not reach the peak performance of DQN in the first case. 642

643 6 Transversal remarks

The contents of previous sections, although presented per application, helps identify trends regarding several technical aspects of state-of-the-art contributions in DRL for fluid flow problems. The present section underlines some of the latter, and raises open questions regarding possible future improvements in the field.

648 6.1 Taking on the challenges

Based on this review, we deem there is a good understanding of the key issues relevant to fluid 649 flow problems. Many of the compiled references are primarily aimed at proving either feasibility 650 651 in such or such sub-domain, or beyond state-of-the-art performance of such or such algorithm, but several milestone contributions assess the ability of novel developments to increase the complexity of 652 the problems presented to the DRL agent. Among the challenges listed in section 3, computational 653 efficiency [40], stochasticity [32, 49, 48] and partial observability [46] have received the most attention, 654 but robustness and delays remain largely ignored (save for the unique combination of stochasticity 655 and post-action time delays examined in [49]), even though real-world environments likely feature all 656 mechanisms in strong interaction one with another. 657

A lot has been achieved in a short period of time, but many related issues remain to be addressed for 658 which the RL literature provides a number of a off-the-shelf methods already proved fruitful in different 659 context (mostly robotics), that could help reach even higher levels of performance and robustness. 660 Typical examples include learning a model of the environment in such a way that errors in the model do 661 not degrade the asymptotic performance [71, 72], or wrapping redundant states into equivalent classes 662 of canonical spaces [73] to increase the data efficiency; using data augmentation and randomization 663 techniques to train over a wide distribution of states [74, 75] or partitioning the initial state distribution 664 and training different policies later to be merged [76] to alleviate stochasticity; optimizing for worst 665 case expected return objectives [77] or pursuing soft-robustness [78] to improve robustness; using the 666 667 frameworks of partially observable Markov decision process [79] and delay-aware Markov Decision Process [80] to account for partial observability and delayed dynamics. 668

Table 4: Usage frequency of different DRL algorithms in the articles considered in the present review. The colours represent the popularity of the method in the domain, a darker colour indicating a more frequent usage. PPO is obviously the most spread method, most probably due to several open-source releases.

DQN	DDQN	A2C	PPO	TRPO	DDPG	TD3	PPO-1/PBO	Others
4	2	1	10	1	3	2	3	3

669 6.2 Providing guidelines for the selection of the DRL algorithm

An obvious preference for policy gradient techniques appears from the review, with PPO the clear-cut 670 go-to algorithm; see table 4. This is noteworthy because PPO is an on-policy algorithm, that updates 671 the policy used to generate the training data (in contrast to off-policy algorithms, that also learn from 672 data generated with other policies). PPO is generally acknowledged to improve the sample efficiency 673 of regular actor-critic techniques, but there could be a fad component to this rise to prominence 674 (partly attributable to the early open-source code release of several projects relying on this technique 675 [14, 40, 43]), given that off-policy methods are expected to have even higher sample efficiency, and that 676 most authors fail to explain the rationale for choosing a particular algorithm over another. Given the 677 high CPU requirements of CFD solvers (that remains an important limitation regarding the application 678 of DRL to 3-D flows of engineering importance), this calls for more careful, consistent and systematic 679 testing of state-of-the-art on- and off-policy techniques in a fluid mechanics context. At the time of 680 writing, only two such comparison studies are available in the literature, namely PPO vs. TD3 in [58], 681 and DQN vs. MFEC in [51]. 682

683 6.3 Fighting the reproducibility crisis

DRL a very fast-moving field, and as the number of contributions is growing, it becomes harder 684 and harder to make a proper comparison between DRL algorithms, all the more so as a bevy of 685 algorithms have been developed, to be used from dedicated libraries (e.g. Tensorforce [62], Stable 686 Baselines [64], OpenAI Baselines [63]) or implemented in-house (which relates to 10 out of the 32 687 reviewed contributions). Compounding the matter are the high amount of time needed to train DRL 688 agents, that creates a high barrier for reevaluation of previous work,; the general lack of complete 689 information regarding the network architecture (e.q. size and depth of the hidden layers, activation690 functions, normalization, initialization) and training procedure (e.g. optimizer, batch size, number of 691 epoch per update, update frequency, learning rate); and (for numerical environments) the additional 692 variance in the numerical solutions themselves. 693

Encouraging the open sourcing of appropriate code on public git repositories is thus a critical step to ensure the reproducibility and durability of the developments, to maximize their impact, and to ultimately help establish DRL as a mature and stable technique for the analysis and design of complex flow systems. In this respect, it is disappointing to note that only 9 out of the 32 studies compiled in this review have come with such open-source releases [14, 41, 43, 44, 45, 40, 16, 59, 15]. Creating and providing exhaustive benchmark datasets and metrics is another alternative that would certainly add value to the community, and lay the ground for solid further developments in the field.

701 6.4 Other research gaps

The present review has also allowed us to identify several other important gaps to consider when evaluating the progress of DRL for practically meaningful fluid mechanics.

Network architecture: almost all provided references use fully-connected networks, with two or three
 hidden layers, each holding a number of neurons in the range from a few tens to a few hundreds.
 Nonetheless, our review did not reveal any large-scale study of the impact of the network architecture
 on the agent performance, and the choice remains most often empirical. The single-step method used
 in [32] is especially interesting in this regards, as it succeeds in learning optimal state-independent
 policies from extremely small networks. It should also be noted that the successful use of LSTM cells
 instead of regular fully-connected networks was advertised in swimming applications [13, 52], and that

additional comparative experiments on different problems could lead to a more systematic use of sucharchitectures.

State space dimensionality: in some cases, state selection seems arbitrary, which can lead to either
(i) incomplete observations or (ii) a too large inputs to the actor, which can be detrimental to learning.
Specific methods have been proposed to tackle this issue, either by adding an intelligent state selection
mechanism [46], or by exploiting state compression [57]. Shall they be pursued further, such efforts
could lead to systematic techniques for state input from CFD environments.

Action space dimensionality: in most contributions, the dimension of the action space remained
 limited, usually between 1 and 3. In this context, Belus *et al.* showed that exploiting the physical
 invariants of the problem was a particularly efficient way to tackle action spaces of larger dimensions
 (up to 20) [15].

• Time granularity: the frequency at which the agent interacts with its environment is usually set based on physical considerations, but the ratio of the typical physical time scale to the action timestep remains highly variable from one contribution to another (even for very similar cases; see table 3). Since this hyper-parameter can dramatically affect the attainable performance of the agent and the difficulty of the learning task (too large intervals lead to inefficient actions, while too small intervals hinder the learning process), the development of systematic selection criteria is another aspect that could benefit the community and help close the gap with real-world testing.

729 7 Conclusion

In the present review, the contributions of the last six years in the field of deep reinforcement learning 730 applied to fluid mechanics problems were presented. The type of application, its complexity, the choice 731 of control methods as well as their associated technical choices were analyzed and compared across 732 the different contributions. Several trends and general rules of thumb currently in use in the domain 733 were pointed out, while unusual choices and techniques were highlighted. This systematic work aims 734 at providing a general frame of the existing usages and techniques to the researchers working in the 735 domain, but also to help newcomers identify standard approaches and state-of-the-art performance 736 level in the field of DRL-based control for fluid dynamics. 737

Overall, impressive performances were observed in multiple complex control tasks. Yet, a large amount 738 of technical questions remain unanswered, and serious efforts remain to be provided by the community 739 in order to efficiently tackle cases of industrial-level complexity within reasonable time. In the pursue 740 of this goal, the access to efficient CFD solvers and to large computational resources remains an issue 741 to many teams. In this perspective, the ability to successfully transfer agents from numerical to 742 experimental environments remains to be explored more thoroughly, as the literature dealing with the 743 coupling of DRL with experimental configurations remains, to this day, extremely scarce. It makes no 744 doubt that the upcoming years will see the mastering of these obstacles, supported by the constant 745 progress made in the DRL field and driven by the numerous industrial challenges that could benefit 746 from it. 747

748 Acknowledgements

This work is supported by the Carnot M.I.N.E.S. Institute through the M.I.N.D.S. project.

750 **References**

- [1] W. Rawat and Z. Wang. Deep convolutional neural networks for image classification: a comprehensive review. Neural Computation, 29:2352–2449, 2017.
- [2] A. Khan, A. Sohail, U. Zahoora, and A. S. Qureshi. A survey of the recent architectures of deep convolutional neural networks. Artificial Intelligence Review, pages 2352–2449, 2020.
- [3] A. B. Nassif, I. Shahin, I. Attili, M. Azzeh, and K. Shaalan. Speech recognition using deep neural networks: a systematic review. IEEE Access, 7:19143–19165, 2019.
- [4] J. Gui, Z. Sun, Y. Wen, D. Tao, and J. Ye. A review on generative adversarial networks: algorithms, theory, and applications. arXiv preprint arXiv:2001.06937, 2020.

- [5] L. Pinto, M. Andrychowicz, P. Welinder, W. Zaremba, and P. Abbeel. Asymmetric actor critic for image-based robot learning. arXiv preprint arXiv:1710.06542, 2017.
- [6] D. Bahdanau, P. Brakel, K. Xu, A. Goyal, R. Lowe, J. Pineau, A. Courville, and Y. Bengio. An actor-critic algorithm for sequence prediction. arXiv preprint arXiv:1607.07086, 2016.
- [7] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller.
 Playing Atari with deep reinforcement learning. arXiv preprint arXiv:1312.5602, 2013.
- [8] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker,
 M. Lai, A. Bolton, Y. Chen, T. Lillicrap, F. Hui, L. Sifre, G. van den Driessche, T. Graepel, and
 D. Hassabis. Mastering the game of Go without human knowledge. Nature, 550, 2017.
- [9] A. Kendall, J. Hawke, D. Janz, P. Mazur, D. Reda, J.-M. Allen, V.-D. Lam, A. Bewley, and
 A. Shah. Learning to drive in a day. arXiv preprint arXiv:1807.00412, 2018.
- [10] A. Bewley, J. Rigley, Y. Liu, J. Hawke, R. Shen, V.-D. Lam, and A. Kendall. Learning to drive from simulation without real world labels. arXiv preprint arXiv:1812.03823, 2018.
- [11] W. Knight. Google just gave control over data center cool-772 ing to an AI. http://www.technologyreview.com/s/611902/ 773 google-just-gave-control-over-data-center-cooling-to-an-ai/, 2018. 774
- [12] G. Novati, S. Verma, D. Alexeev, D. Rossinelli, W. M. van Rees, and P. Koumoutsakos. Syn chronisation through learning for two self-propelled swimmers. <u>Bioinspiration & Biomimetics</u>, 12(3):036001, 2017.
- [13] S. Verma, G. Novati, and P. Koumoutsakos. Efficient collective swimming by harnessing vortices through deep reinforcement learning. <u>Proceedings of the National Academy of Sciences</u>, 115(23):5849–5854, 2018.
- [14] J. Rabault, M. Kuchta, A. Jensen, U. Réglade, and N. Cerardi. Artificial neural networks trained
 through deep reinforcement learning discover control strategies for active flow control. Journal of
 Fluid Mechanics, 865:281–302, 2019.
- [15] V. Belus, J. Rabault, J. Viquerat, Z. Che, E. Hachem, and U. Reglade. Exploiting locality and translational invariance to design effective deep reinforcement learning control of the 1-dimensional unstable falling liquid film. AIP Advances, 9(12):125014, 2019.
- [16] J. Viquerat, J. Rabault, A. Kuhnle, H. Ghraieb, A. Larcher, and E. Hachem. Direct shape
 optimization through deep reinforcement learning. Journal of Computational Physics, 428:110080,
 2021.
- [17] P. Garnier, J. Viquerat, J. Rabault, A. Larcher, A. Kuhnle, and E. Hachem. A review on deep
 reinforcement learning for fluid mechanics. Computers & Fluids, 225:104973, 2021.
- [18] W. Zhang J. Rabault, F. Ren. Deep reinforcement learning in fluid mechanics: A promising method for both active flow control and shape optimization. Journal of Hydrodynamics, 32:234– 246, 2020.
- [19] R. S. Sutton and A. G. Barto. <u>Reinforcement Learning: An Introduction</u>. MIT Press, Cambridge,
 MA, 2018.
- R. Bellman and S. E. Dreyfus. <u>Applied dynamic programming</u>. Princeton University Press Princeton, N.J, 1962.
- R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3):229-256, 1992.
- [22] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. Neural Networks, 2(5):359–366, 1989.
- ⁸⁰³ [23] H. T. Siegelmann and E. D. Sontag. On the computational power of neural nets. Journal of ⁸⁰⁴ Computer and System Sciences, 50(1):132–150, 1995.
- ⁸⁰⁵ [24] I. Goodfellow, Y. Bengio, and A. Courville. The Deep Learning Book. MIT Press, 2017.
- [25] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Ried miller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King,
 D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforce-
- 809 ment learning. <u>Nature</u>, 518, 2015.

- [26] T. Schaul, J. Quan, I. Antonoglou, and D. Silver. Prioritized experience replay. <u>arXiv preprint</u> arXiv:1511.05952, 2016.
- [27] H. van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double q-learning. <u>arXiv</u>
 preprint arXiv:1509.06461, 2015.
- [28] V. Mnih, A. Puigdomènech Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver,
 and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. arXiv preprint
 arXiv:1602.01783, 2016.
- [29] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel. Trust region policy optimization.
 arXiv preprint arXiv:1502.05477, 2015.
- [30] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
 algorithms. arXiv preprint arXiv:1707.06347, 2017.
- [31] S.Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. arXiv preprint arXiv:1802.09477, 2018.
- [32] H. Ghraieb, J. Viquerat, A. Larcher, P. Meliga, and E. Hachem. Single-step deep reinforcement
 learning for open-loop control of laminar and turbulent flows. <u>arXiv preprint arXiv:2006.02979</u>,
 2020.
- [33] E. Hachem, H. Ghraieb, J. Viquerat, A. Larcher, and P. Meliga. Deep reinforcement learning
 for the control of conjugate heat transfer with application to workpiece cooling. arXiv preprint
 arXiv:2011.15035, 2020.
- ⁸²⁹ [34] N. Hansen. The cma evolution strategy: a rtutorial. arXiv preprint arXiv:1604.00772, 2016.
- [35] J. Viquerat, R. Duvigneau, P. Meliga, A. Kuhnle, and E. Hachem. Policy-based optimization:
 single-step policy gradient method seen as an evolution strategy. arXiv preprint arXiv:2104.06175,
 2021.
- [36] G. Dulac-Arnold, D. Mankowitz, and T. Hester. Challenges of real-world reinforcement learning.
 arXiv preprint arXiv:1904.12901, 2019.
- [37] G. Dulac-Arnold, N. Levine, D. J. Mankowitz, J. Li, C. Paduraru, S. Gowal, and T. Hester.
 Challenges of real-world reinforcement learning: definitions, benchmarks and analysis. <u>Machine</u>
 Learning, pages 1–50, 2021.
- [38] J. J. Garau-Luis, E. Crawley, and B. Cameron. Evaluating the progress of deep reinforcement
 learning in the real world: aligning domain-agnostic and domain-specific research. arXiv preprint
 arXiv:2107.03015, 2021.
- [39] H. Koizumi, S. Tsutsumi, and E. Shima. Feedback control of karman vortex shedding from a cylinder using deep reinforcement learning. 2018 Flow Control Conference, 2018.
- [40] J. Rabault and A. Kuhnle. Accelerating deep reinforcement learning strategies of flow control
 through a multi-environment approach. Physics of Fluids, 31(9):094105, 2019.
- [41] M. Tokarev, E. Palkin, and R. Mullyadzhanov. Deep reinforcement learning control of cylinder
 flow using rotary oscillations at low reynolds number. Energies, 13(22), 2020.
- [42] H. Xu, W. Zhang, J. Deng, and J. Rabault. Active flow control with rotating cylinders by an artificial neural network trained by deep reinforcement learning. <u>Journal of Hydrodynamics</u>, 32:254–258, 2020.
- [43] H. Tang, J. Rabault, A. Kuhnle, Y. Wang, and T. Wang. Robust active flow control over a range of reynolds numbers using an artificial neural network trained through deep reinforcement learning. Physics of Fluids, 32(5):053605, 2020.
- [44] M. A. Elhawary. Deep reinforcement learning for active flow control around a circular cylinder
 using unsteady-mode plasma actuators. arXiv preprint arXiv:2012.10165, 2020.
- [45] M. Holm. Using deep reinforcement learning for active flow control. Master's thesis, University
 of Oslo, 2020.
- [46] R. Paris, S. Beneddine, and J. Dandois. Robust flow control and optimal sensor placement using
 deep reinforcement learning. arXiv preprint arXiv:2006.11005, 2020.
- [47] S. Qin, S. Wang, and G. Sun. An application of data driven reward of deep reinforcement learning
 by dynamic mode decomposition in active flow control. arXiv preprint arXiv:2106.06176, 2021.

- [48] F. Ren, J. Rabault, and H. Tang. Applying deep reinforcement learning to active flow control in
 weakly turbulent conditions. Physics of Fluids, 33(3):037121, 2021.
- [49] G. Beintema, A. Corbetta, L. Biferale, and F. Toschi. Controlling rayleigh-bénard convection via
 reinforcement learning. Journal of Turbulence, 21(9-10):585-605, 2020.
- [50] X. Y. Lee, A. Balu, D. Stoecklein, B. Ganapathysubramanian, and S. Sarkar. A case study of
 deep reinforcement learning for engineering design: application to microfluidic devices for flow
 sculpting. Journal of Mechanical Design, 141(11), 2019.
- [51] O. J. Dressler, P. D. Howes, J. Choo, and A. J. deMello. Reinforcement learning for dynamic microfluidic control. ACS Omega, 3(8):10084–10091, 2018.
- Y. Zhu, F.-B. Tian, J. Young, J. C. Liao, and J. C. S. Lai. A numerical study of fish adaption
 behaviors in complex environments with a deep reinforcement learning and immersed boundary-lattice boltzmann method. Nature Scientific Reports, 11(1691), 2021.
- [53] L. Yan, X. Chang, R. Tian, N. Wang, L. Zhang, and W. Liu. A numerical simulation method for bionic fish self-propelled swimming under control based on deep reinforcement learning.
 Proceedings of the Institution of Mechanical Engineers, Part C: Journal of Mechanical Engineering Science, 234(17):3397–3415, 2020.
- [54] X. Yan, J. Zhu, M. Kuang, and X. Wang. Aerodynamic shape optimization using a novel optimizer
 based on machine learning techniques. Aerospace Science and Technology, 86:826–835, 2019.
- [55] R. Li, Y. Zhang, and H. Chen. Learning the aerodynamic design of supercritical airfoils through
 deep reinforcement learning. arXiv preprint arXiv:2010.03651, 2020.
- [56] S. Qin, S. Wang, L. Wang, C. Wang, G. Sun, and Y. Zhong. Multi-objective optimization of
 cascade blade profile based on reinforcement learning. Applied Sciences, 11(1), 2021.
- [57] P. Ma, Y. Tian, Z. Pan, B. Ren, and D. Manocha. Fluid directed rigid body control using deep
 reinforcement learning. ACM Transactions on Graphics, 37(4), 2018.
- Y. Xie and X. Zhao. Sloshing suppression with active controlled baffles through deep reinforcement
 learning-expert demonstrations-behavior cloning process. Physics of Fluids, 33(1):017115, 2021.
- [59] G. Novati, H. L. de Laroussilhe, and P. Koumoutsakos. Automating turbulence modelling by
 multi-agent reinforcement learning. Nature Machine Intelligence, 3:87–96, 2021.
- [60] D. Fan, L. Yang, Z. Wang, M. S. Triantafyllou, and G. E. Karniadakis. Reinforcement learning
 for bluff body active flow control in experiments and simulations. <u>Proceedings of the National</u>
 Academy of Sciences, 117(42):26091–26098, 2020.
- [61] S. Shimomura, S. Sekimoto, A. Oyama, K. Fujii, and H. Nishida. Closed-loop flow separation
 control using the deep q-network over airfoil. AIAA Journal, 58(10):4260–4270, 2020.
- [62] A. Kuhnle, M. Schaarschmidt, and K. Fricke. Tensorforce: a tensorflow library for applied reinforcement learning. https://github.com/tensorforce/tensorforce, 2017.
- [63] P. Dhariwal, C. Hesse, O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor,
 Y. Wu, and P. Zhokhov. Openai baselines. https://github.com/openai/baselines, 2017.
- [64] A. Hill, A. Raffin, M. Ernestus, A. Gleave, A. Kanervisto, R. Traore, P. Dhariwal, C. Hesse,
 O. Klimov, A. Nichol, M. Plappert, A. Radford, J. Schulman, S. Sidor, and Y. Wu. Stable
 baselines. https://github.com/hill-a/stable-baselines, 2018.
- [65] M. S. Alnæs, J. Blechta, J. Hake, A. Johansson, B. Kehlet, A. Logg, C. Richardson, J. Ring,
 M. E. Rognes, and G. N. Wells. The fenics project version 1.5. <u>Archive of Numerical Software</u>,
 3(100), 2015.
- [66] A. Lampton, A. Niksch, and J. Valasek. Morphing airfoils with four morphing parameters. In
 AIAA Guidance, Navigation and Control Conference and Exhibit, 2008.
- [67] M. Gazzola, B. Hejazialhosseini, and P. Koumoutsakos. Reinforcement learning and wavelet adapted vortex methods for simulations of self-propelled swimmers. <u>SIAM Journal of Scientific</u> Computing, 36:622–639, 2014.
- [68] M. Gazzola, A. A. Tchieu, D. Alexeev, A. de Brauer, and P. Koumoutsakos. Learning to school
 in the presence of hydrodynamic interactions. Journal of Fluid Mechanics, 789:726–749, 2016.

- [69] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin,
 P. Abbeel, and W. Zaremba. Hindsight experience replay. arXiv preprint arXiv:1707.01495, 2018.
- [70] C. Blundell, B. Uria, A. Pritzel, Y. Li, A. Ruderman, Joel Z Leibo, Jack Rae, Daan Wierstra, and Demis Hassabis. Model-free episodic control. arXiv preprint arXiv:1606.04460, 2016.
- [71] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. arXiv preprint arXiv:1805.12114, 2018.
- [72] J. Buckman, D; Hafner, G. Tucker, E. Brevdo, and H. Lee. Sample-efficient reinforcement learning
 with stochastic ensemble value expansion. arXiv preprint arXiv:1807.01675, 2018.
- [73] C. Wu, A. Kreidieh, E. Vinitsky, and A. M. Bayen. Emergent behaviors in mixed-autonomy
 traffic. In Conference on Robot Learning, pages 398–407, 2017.
- [74] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization
 for transferring deep neural networks from simulation to the real world. In <u>2017 IEEE/RSJ</u>
 international conference on intelligent robots and systems (IROS), pages 23–30, 2017.
- [75] K. Lee, K. Lee, J. Shin, and H. Lee. Network randomization: A simple technique for generalization in deep reinforcement learning. arXiv preprint arXiv:1910.05396, 2019.
- [76] D. Ghosh, A. Singh, A. Rajeswaran, V. Kumar, and S. Levine. Divide-and-conquer reinforcement learning. arXiv preprint arXiv:1711.09874, 2017.
- [77] D. J. Mankowitz, N. Levine, R. Jeong, Y. Shi, J. Kay, A. Abdolmaleki, J. T. Springenberg,
 T. Mann, T. Hester, and M. Riedmiller. Robust reinforcement learning for continuous control
 with model misspecification. arXiv preprint arXiv:1906.07516, 2019.
- [78] E. Derman, D. J. Mankowitz, T. A. Mann, and S. Mannor. Soft-robust actor-critic policy-gradient.
 arXiv preprint arXiv:1803.04848, 2018.
- [79] A. R. Cassandra, L. P. Kaelbling, and M. L. Littman. Acting optimally in partially observable
 stochastic domains. In AAAI, volume 94, pages 1023–1028, 1994.
- [80] B. Chen, M. Xu, L. Li, and D. Zhao. Delay-aware model-based reinforcement learning for continuous control. Neurocomputing, 450:119–128, 2021.