# Migration of self-propelling agent in a turbulent environment with minimal energy consumption

Ao Xu,[1, 2, a)] Hua-Lin Wu,[1] and Heng-Dong Xi[1, 2]

[1)]*School of Aeronautics, Northwestern Polytechnical University, Xi'an 710072, China*

[2)]*Institute of Extreme Mechanics, Northwestern Polytechnical University, Xi'an 710072, China*

(Dated: 17 March 2022)

We present a numerical study of training a self-propelling agent to migrate in the unsteady flow environment. We control the agent to utilize the background flow structure by adopting the reinforcement learning algorithm to minimize energy consumption. We considered the agent migrating in two types of flows: one is simple periodical double-gyre flow as a proof-of-concept example, while the other is complex turbulent Rayleigh-Bénard convection as a paradigm for migrating in the convective atmosphere or the ocean. The results show that the smart agent in both flows can learn to migrate from one position to another while utilizing background flow currents as much as possible to minimize the energy consumption, which is evident by comparing the smart agent with a naive agent that moves straight from the origin to the destination. In addition, we found that compared to the double-gyre flow, the flow field in the turbulent Rayleigh-Bénard convection exhibits more substantial fluctuations, and the training agent is more likely to explore different migration strategies; thus, the training process is more difficult to converge. Nevertheless, we can still identify an energy-efficient trajectory that corresponds to the strategy with the highest reward received by the agent. These results have important implications for many migration problems such as unmanned aerial vehicles flying in a turbulent convective environment, where planning energy-efficient trajectories are often involved. [a]

[a)]Electronic mail: Corresponding author: axu@nwpu.edu.cn (Ao Xu)

# I.  INTRODUCTION

Soaring birds and gliders often use spatially and temporally localized warm rising atmospheric currents to stay aloft and fly higher.[1] The so-called *thermal soaring* behavior[2] can save a vast amount of energy by minimizing the flapping of wings for birds or motor power supplies for gliders over long distances.[3] For example, the Andean condor can even soar over 5h (covering about 172 km) without flapping.[4] However, the air currents in the troposphere are turbulent; thus, using wobbly gusts of air to stay airborne has not always been a simple task. Ideal conditions for thermal soaring typically occur when a strong temperature gradient between the surface of the Earth and the top of the atmospheric boundary layer creates convective thermals. The convective thermals exhibit general turbulence characterized by strongly fluctuating flow velocities.

Learning about details of thermal soaring can improve understanding of the main features of flight trajectories and optimization strategies. In 1958, MacCready[5] proposed a theory on flight optimization and gave a gliding polar curve (the relationship between horizontal speed and the corresponding vertical one) to calculate the best slope to take before an upcoming thermal. Since then, gliders have tried to adjust their gliding speed to the expected thermal climb rate according to their polar curve. With the aid of measured polar curves, Akos *et al.*[6] found that there are relevant common features in the way that falcons and the world's leading paraglider pilots use thermals, which are also close to the optimal soring strategy predicted by MacCready's theory. To apply the above soaring strategy for an unmanned aerial vehicle (UAV) to take advantage of thermals, Allen and Lin[7] adopted autonomous soaring algorithms to detect and exploit thermals. He used the aircraft's total energy state to detect and soar within thermals and the estimated thermal size and position to calculate guidance commands for soaring flight. On the other hand, reinforcement learning (RL) methods are promising to deliver effective strategies of soaring flight. For example, Wharington and Herszberg[8] used a neural-based algorithm to locate the thermal core. However, they only considered the learning problem of finding the center of a stationary thermal without turbulence. Later, Akos *et al.*[9] demonstrated that such simple rules would fail in the presence of velocity fluctuations. Thus, soaring strategies that could work in real turbulent work are urgently needed. To identify effective soaring strategies of flight in turbulent flows, Reddy *et al.*[10,11] used reinforcement learning algorithms to train gliders to travel through complex choppy air currents. Environment cues, such as an increase in the twisting force of the wind that indicates rising air, can be sensed. However, in the work of Reddy *et al.*, their motivation is to train the glider to employ

spiraling patterns to ascend higher in regions of strong upwelling currents. In practical flight tasks, we would also expect the UAV to fly from one position to another while utilizing thermal soaring as much as possible to extend the flight duration and reduce the energy consumption.[12]

Previously, reinforcement learning algorithms have also been applied to navigate micro gravitactic swimmers to escape local fluid traps and reach the highest altitude,[13,14] to accumulate in regions of intense negative vorticity,[15] and to reach a target position with minimal time.[16–18] Shape effects (e.g., elliptical or ellipsoidal shape) of asymmetric swimmers have also been considered with the goal of moving upwards.[19,20] In this work, our motivation is to train an active self-propelling agent to migrate from one position to another with minimal energy cost. This motivation stems from reducing the energy consumption for UAVs during their flight. The rest of this paper is organized as follows. In Sec. II, we present details for the reinforcement learning algorithm to train the self-propelling agent to find an energy-efficient trajectory. In Sec. III, we first train the agent in the unsteady double-gyre flow,[21] which is a simple periodic flow environment with an analytical solution for the flow field, as a proof-of-concept example. In Sec. IV, we then train the agent in the turbulent Rayleigh-Bénard convection,[22] which is a much more complex turbulent environment with a strongly fluctuating flow field as a paradigm for migration in the atmosphere or the ocean. In Sec. V, the main findings of this work are summarized.

## II. DYNAMICS OF THE SELF-PROPELLING AGENT

### A. Optimal control via the reinforcement learning algorithm

We aim to plan an energy-efficient trajectory for the self-propelling agent in the unsteady flow environment. Traditional approaches, such as the optimal navigation theory,[23,24] may be sensitive to small disturbances in the chaotic system. Thus, we adopt the emerging reinforcement learning (RL) algorithm to optimize the trajectory in this work.[25–28] In the RL algorithm, the agent observes the state of the environment and then decides on an action to take. If the agent receives a reward (or a penalty) for that action, it is more likely to repeat (or forego) the action in the future. Overall, the agent learns by trial and error and eventually achieves its goal.[29,30] In the double-gyre flow, the observation variables include the background flow velocity $\mathbf{u}_{\text{fluid}}$, the agent's spatial coordinates $\mathbf{x}_{\text{agent}}$, and the current time $t$; in the turbulent Rayleigh-Bénard convection, we also include the fluid temperature $T$ in addition to the above-mentioned observation variables. The action variable

is that the agent generates propelling velocity of $\mathbf{u}_{\mathrm{propel}}$.

The model-free reinforcement learning algorithms can generally be classified into two categories: one is the policy optimization method, in which the parameter $\theta$ is optimized to maximize the performance objective $J(\pi_\theta)$, and the other is the Q-learning method, in which the agent takes action $a$ that tried to maximize the optimal action-value function, i.e., $a(s) = \arg\max_a Q_\theta(s,a)$. Here, $s$ denotes the state of the environment, $\pi_\theta$ denotes the parameterized stochastic policy, and $Q_\theta(s,a)$ approximates the optimal action-value function $Q^*(s,a)$. However, the policy optimization method is inefficient in sampling, because it cannot reuse data to train the model, while the Q-learning method tends to be less stable because it indirectly optimizes the agent performance.[31,32] A trade-off between these two methods is the soft actor-critic (SAC) method,[33] in which the actor aims to maximize the expected reward (i.e., succeed at the task) while also maximizing entropy (i.e., acting as randomly as possible). In entropy-regularized reinforcement learning, the optimization problem can be described as

$$\pi^*(\theta) = \arg\max_\pi \underset{\tau \sim \pi}{E} \left[ \sum_{t=0}^{\infty} \left( R(s_t, a_t, s_{t+1}) + \alpha H(\pi(\cdot|s_t))) \right) \right] \tag{1}$$

In the above equation, $\pi^*$ is the optimal policy. The reward function $R$ depends on the current state of the environment $s_t$, the action just taken $a_t$, and the next state of the environment $s_{t+1}$. $\alpha$ is the trade-off coefficient. The entropy $H$ of $\tau$ is computed from its distribution $\pi$ as $H(\pi(\cdot|s_t)) = \underset{\tau \sim \pi}{E} \left[ -\log \pi(\tau) \right]$. More details on the SAC method can be found in Ref.[33].

In this work, we assume the rewards gained by the agent are simultaneously affected by its current state, energy consumption, and time consumption. We design the reward function as

$$r(t) = r_s(t) + r_e(t) + r_h(t) \tag{2}$$

Here, $r_s$ denotes the reward contributed by the current state of the agent. We assume that if the agent migrates out of the flow domain, it will receive a penalty of -10; if the agent is getting closer to the destination, it will receive a basic reward of $e_{\mathrm{basic}}$. Thus, we can express $r_s$ as

$$r_s = \begin{cases} -10, & \text{agent is out of the flow domain} \\ e_{\mathrm{basic}}, & \left\| \mathbf{x}_{\mathrm{agent}}^{t+1} - \mathbf{x}_{\mathrm{goal}} \right\|_2 < \left\| \mathbf{x}_{\mathrm{agent}}^{t} - \mathbf{x}_{\mathrm{goal}} \right\|_2 \\ 0, & \text{otherwise} \end{cases} \tag{3}$$

Here, $r_e$ denotes the reward contributed by the energy consumption of the agent. We assume that if the propelling velocity of the agent $\mathbf{u}_{\mathrm{propel}}$ is in alignment with that of the background flow

$\mathbf{u}_{\text{fluid}}$, namely, the angle between these two vectors is $\alpha \leq 90°$, the agent will receive a reward of $e_{\text{basic}} + (e_{\text{max}} - e)$, where $e = 0.5 \left\| \mathbf{u}_{\text{propel}} \right\|^2$ and $e_{\text{basic}} = e_{\text{max}} = 0.5(\left\| \mathbf{u}_{\text{propel}} \right\|)^2_{\text{max}}$, suggesting that when the agent migrates, it follows the background flow direction, the higher the agent generates propelling velocity, the lower the reward it receives; otherwise, it will receive a penalty of $-(e_{\text{basic}} + e)$, suggesting that when the agent migrates against the background flow direction, the higher the agent generates propelling velocity, the higher the penalty it receives. Thus, we can express $r_e$ as

$$r_e = \begin{cases} e_{\text{basic}} + (e_{\text{max}} - e), & \text{for } 0° \leq \alpha \leq 90° \\ -(e_{\text{basic}} + e), & \text{for } 90° < \alpha \leq 180° \end{cases} \tag{4}$$

Here, $r_h$ denotes the reward contributed by the time consumption of the agent. We assume that if the agent cannot reach the destination within a maximum time of $t_{\text{max}}$, it will receive a penalty of -5; if the agent is within $\delta_0$ from the destination (here, $\delta_0$ denotes a small threshold value), we assume the agent reaches the destination and it will receive a reward that is inversely proportional to the time taken (with coefficient $\varepsilon$) during the migration, suggesting the sooner the agent reaches the destination, the higher the reward it receives. Thus, we can express $h_t$ as

$$r_h = \begin{cases} -5, & t \geq t_{\text{max}} \\ \varepsilon(t_{\text{max}} - t), & \left\| \mathbf{x}^t_{\text{agent}} - \mathbf{x}_{\text{goal}} \right\|_2 < \delta_0 \\ 0, & \text{otherwise} \end{cases} \tag{5}$$

Our specially designed reward function also implies that when the agent migrates toward the destination (i.e., $\left\| \mathbf{x}^{t+1}_{\text{agent}} - \mathbf{x}_{\text{goal}} \right\|_2 < \left\| \mathbf{x}^t_{\text{agent}} - \mathbf{x}_{\text{goal}} \right\|_2$) following the background flow direction (i.e., $0° \leq \alpha \leq 90°$), in each timestep, reducing energy consumption will be the agent's primary objective, while approaching the destination will be its second objective [because $e_{\text{basic}} + (e_{\text{max}} - e) > e_{\text{basic}}$]. In addition, if the agent has to reach the destination within a short time (i.e., $t_{\text{max}}$ is not long enough for the agent to freely explore the environment), in each episode, approaching the destination will be the agent's primary objective, while reducing energy consumption will be its second objective [because $E\left[\sum r_h\right] > E\left[\sum (r_s + r_e)\right]$].

## B. Kinematic model for the self-propelling agent

We restricted the maximum propelling velocity of the agent $\mathbf{u}_{\text{propel}}$ to be less than the largest background flow velocity, such that intelligent planning of the agent can well utilize the background flow structure, also to mimic the limited propulsion available in real-world scenarios. We

assume that, without control, the velocity of the agent equals the velocity of the background fluid flow $\mathbf{u}_{\text{fluid}}$; meanwhile, the agent can take action to generate its own relative velocity $\mathbf{u}_{\text{propel}}$. Then, with control, we can model the agent's velocity in the unsteady flow as $\mathbf{u}_{\text{agent}} = \mathbf{u}_{\text{fluid}} + \mathbf{u}_{\text{propel}}$. The position of the agent is updated via the relation $d\mathbf{x}_{\text{agent}}/dt = \mathbf{u}_{\text{agent}}$. It is worth mentioning that the present kinematic model for the agent is far from a realistic one in the industry. Here, we adopt this simple kinematic model to disentangle the coupling between the chaotic flow of the carrier fluid and the complex motion of the agent. On the other hand, a more complex kinematic model for the self-propelling agent that includes inertial and rotational dynamics,[34] flapping motion,[35,36] or even the flexible motion of the propelling agent[37–39] can be considered in the future work.

## III.   MIGRATION IN THE UNSTEADY DOUBLE-GYRE FLOW

### A.   Numerical simulation of unsteady double-gyre flow

The double-gyre flow field has the analytical description of the velocity field. It has been used to study mixing and coherent structures in large-scale ocean circulation. The flow is defined on a nondimensional domain of $[0,2] \times [0,1]$. The double-gyre velocity field is derived from the stream function

$$\phi(x,y,t) = A \sin\left[\pi f(x,t)\right] \sin(\pi y) \tag{6}$$

and the resulting velocity field is

$$u(x,y,t) = -\pi A \sin\left[\pi f(x,t)\right] \cos(\pi y) \tag{7a}$$

$$v(x,y,t) = -\pi A \cos\left[\pi f(x,t)\right] \sin(\pi y) \tag{7b}$$

In the above equations, the time dependency is introduced by

$$f(x,t) = a(t)x^2 + b(t)x \tag{8}$$

with time-dependent coefficient

$$a(t) = \varepsilon \sin(\omega t), \quad b(t) = 1 - 2\varepsilon \sin(\omega t) \tag{9}$$

Here, $A$ determines the magnitude of the velocity vectors, $\varepsilon$ is the amplitude of the motion of the separation point on the $x$-axis and $\omega$ is the angular oscillation frequency. Unless otherwise mentioned, we adopt the parameter sets of $A = 0.1$, $\varepsilon = 0.25$, and $\omega = 2\pi/10$ as default values.[21]

## B. Training results and discussion

In training, we restrict the maximum propelling velocity of the agent along either $x$- or $y$-direction to be less than $A = 0.1$, which leads to $\|\mathbf{u}_{\text{propel}}\|^2 \leq 0.02$. We show the instantaneous trajectories for the smart agent in the double-gyre flow in Fig. 1 (Multimedia view). We released the agent at (2.0, 1.0) position, namely the top-right corner in the domain (marked by the blue square in the plot). The agent's goal is to reach the destination position of (0.25, 0.8), marked by the red star in the plot, with minimal energy consumption. Here, we assume the agent reaches the destination so long as its position is within $\delta_0 = 0.05$ from (0.25, 0.8). Initially, the smart agent will move leftward to utilize the horizontal currents in the top-right region [see Fig. 1(a)]. When the smart agent reaches the top-left corner of the right gyre at the position around (1, 0.9), it will move downward to utilize the vertical currents [see Fig. 1(b)]. After that, it will drift into the bottom-right corner of the left gyre at the position around (1, 0.1) [see Fig. 1(c)] and then follow the horizontal leftward and vertical upward current to reach the destination [see Fig. 1(d)]. The smart agent follows background currents as much as possible, except at the ridge between the two gyres around $x = 1$, where the agent will consume more energy and drift from the right gyre into the left gyre. Because the structure of the flow leaves its mark on the trajectories of the agents carried by turbulent flows,[40] to quantitatively describe the relationship between the orientation of the propelling velocity vector (i.e., $\theta_{\text{propel}}$) and the orientation of the fluid velocity vector (i.e., $\theta_{\text{fluid}}$), we calculate their cross-correlation coefficient as

$$C = \langle [\theta_{\text{propel}}(t) - \langle \theta_{\text{propel}} \rangle] [\theta_{\text{fluid}}(t) - \langle \theta_{\text{fluid}} \rangle] \rangle / (\sigma_{\theta_{\text{propel}}} \sigma_{\theta_{\text{fluid}}}) \tag{10}$$

The resulting correlation coefficient of 0.60 suggests that the orientations of the propelling velocity vector and fluid velocity vector are statistically relevant. We also visualize the propelling velocity vector and fluid velocity vector of the smart agent in the Appendix.

To demonstrate that the smart agent can, indeed, reduce energy consumption, we compare the strategy of the smart agent with the naive agent. Here, the naive agent moves straight from the origin to the destination with constant velocity, which may be the "simplest" way for an agent to migrate from one position to another. To make a fair comparison on the energy consumption, we set the naive agent spending the same amount of time $t_{\text{total}}$ as the smart agent migrating from the origin to the destination. Thus, the velocity magnitude of the agent is $\|\mathbf{u}_{\text{agent}}\| = \|\mathbf{x}_{\text{goal}} - \mathbf{x}_{\text{start}}\|/t_{\text{total}}$ and its direction point from the origin to the destination. In Fig. 2, we show the trajectories of the naive agent and the smart agent, which are color-coded by the instantaneous
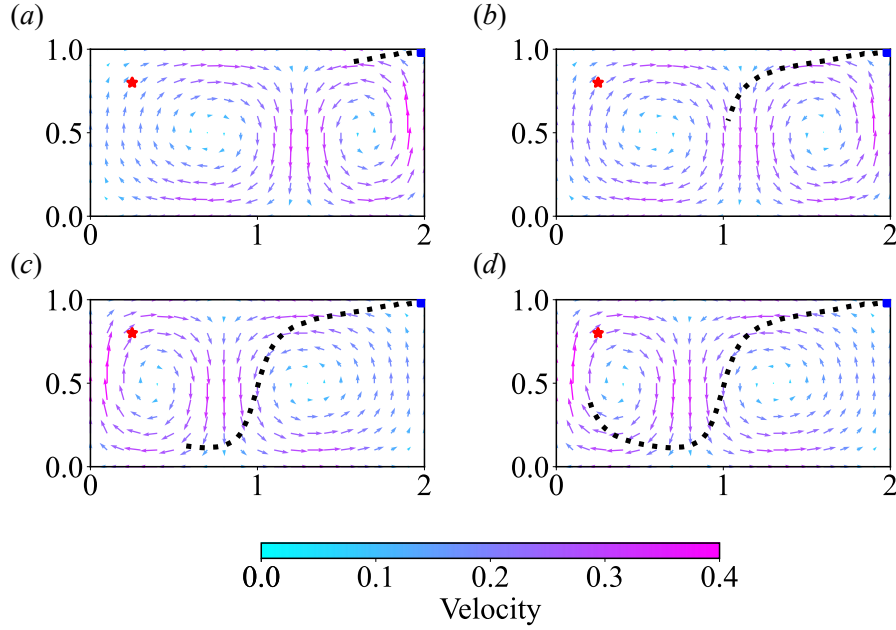
FIG. 1. Trajectory (black dotted line) of the smart agent in the double-gyre flow at $(a)$ $t = 2.2$, $(b)$ $t = 4.2$, $(c)$ $t = 6.7$, $(d)$ $t = 8.2$. The vectors denote the velocity field of the double-gyre flow, and they are color-coded by the velocity magnitude. (Multimedia view)
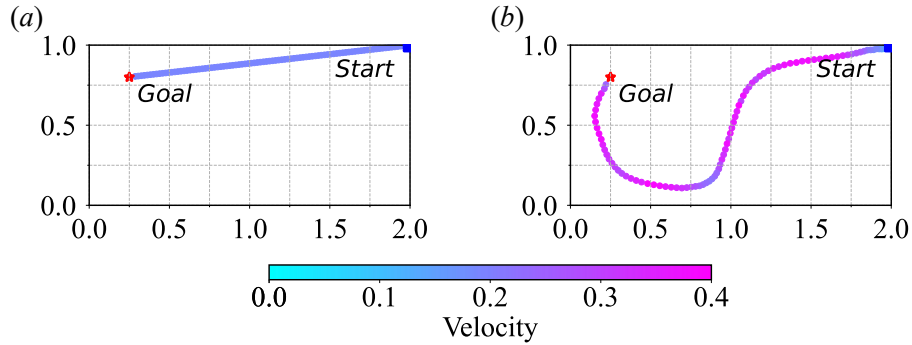


FIG. 2. Comparison of trajectories in the double-gyre flow for $(a)$ a naive agent moving straightly and $(b)$ a smart agent utilizing the flow structure to save energy. The trajectories are color-coded by the instantaneous velocity magnitude.

velocity magnitude. We can see from Fig. 2 that the naive agent migrates slower than the smart agent, because the naive agent's travel distance is shorter than that of the smart agent. Although the smart agent migrates faster, it does not indicate that the smart agent will consume more energy because the smart agent can utilize the flow currents to save energy.

To make a quantitative comparison, we plot the time series of accumulative energy consumed
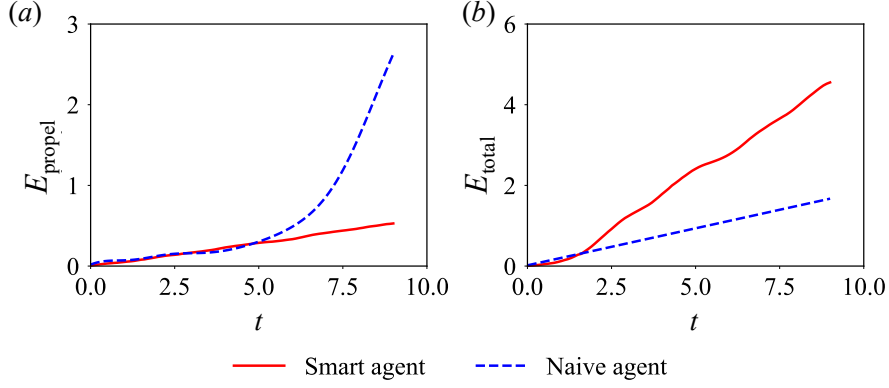
FIG. 3. Comparison of the accumulative (*a*) energy consumption $E_{\text{propel}}$ and (*b*) total kinetic energy $E_{\text{total}}$ of the smart agent and the naive agent in the double-gyre flow.

by the agents, which is calculated as

$$E_{\text{propel}}(t) = \int_0^t \frac{1}{2} \left\| \mathbf{u}_{\text{propel}}(\tau) \right\|^2 d\tau \tag{11}$$

In this work, we assume the agent generates its own velocity of $\mathbf{u}_{\text{propel}}$, which is responsible for its energy consumption. As shown in Fig. 3(a), we can see that both agents consume almost the same amount of energy during the initial period (i.e., around $t < 6$), which is due to similar trajectories in the flow field (see Fig. 2), and both agents migrate in the top area of the right gyre. After that (i.e., around $t > 6$), the smart agent moves downward to utilize the flow currents, while the naive agent continues to move straight toward the destination. Crossing the ridge of the gyre and migrating reversely against the flow currents will require substantial energy consumption, as evident from the much higher energy consumption at $t < 6$ for the naive agent [see Fig. 3(a)]. We also compare the accumulative total kinetic energy of the agents [see Fig. 3(b)], which is calculated as

$$E_{\text{total}}(t) = \int_0^t \frac{1}{2} \left\| \mathbf{u}_{\text{agent}}(\tau) \right\|^2 d\tau \tag{12}$$

After reaching the destination, the accumulative total kinetic energy of the smart agent is more than twice that of the naive agent, while the smart agent only consumed almost one-fifth of the energy, suggesting the smart agent can efficiently utilize the energy provided by the background flow.

In the following, we test the robustness of the energy-efficient strategies with respect to the flow control parameters. We varied $\varepsilon$ in the range of $0.025 \leq \varepsilon \leq 2.5$ while keeping $A$ and $\omega$ fixed as the

9

default values. The resulting optimized trajectories show minor differences, suggesting energy-efficient strategy is robust with the changes in the magnitude of oscillation in the $x$-direction. Similarly, we varied $\omega$ in the range of $0.2\pi/10 \leq \omega \leq 2\pi$ while keeping $A$ and $\varepsilon$ fixed as the default values. The resulting optimized trajectories show minor differences, suggesting energy-efficient strategy is also robust with the changes in the angular oscillation frequency. We then vary $A$ in the range of $0.01 \leq A \leq 1$ while keeping $\varepsilon$ and $\omega$ fixed as the default values. As shown in Fig. 4, the resulting optimized trajectories show significant differences: when the carrier fluid flows at a low speed [see Figs. 4(a) and 4(b), for $A = 0.01$ and $A = 0.05$, respectively], the migration of the agent mostly depends on its own propulsion, and it slowly approaches the destination or even cannot reach the destination; when the carrier fluid flows at a high speed [see Figs. 4(e) and 4(f), for $A = 0.2$ and $A = 1.0$, respectively], the migration of the agent is significantly influenced by the carrier flow, and it may miss the destination or even drift out the flow domain. We also slightly increase or decrease 20% of the $A$ value [see Figs. 4(c) and 4(d), for $A = 0.08$ and $A = 0.12$, respectively], and the resulting optimized trajectories show minor differences. Thus, the optimized energy-efficient strategy is not sensitive to small perturbation of the magnitude of the velocity vectors, but it certainly changes when the magnitude of the velocity vectors varies in a larger range.

We finally show the episode rewards as a function of time steps during the training process. In Fig. 5, the discrete blue dots represent the accumulated rewards obtained by the agents at each episode, and the orange line represents the smooth average of the rewards during a short-time window of 100 episodes. A reward value around -10 indicates that the agent migrates outside of the flow domain received a penalty; a reward value between -5 and 0 indicates that the agent neither migrates outside of the flow domain nor reaches the destination, but it is trapped in a gyre and its trajectories form loops; a reward value around 5 indicates that the agent reaches the destination, yet, it takes a high energy consumption and a long migration time; a reward value around 10 indicates that the agent can successfully reach the goal with minimal energy cost and migration time. We can see from Fig. 5 that, initially (i.e., for timesteps less than 50 000), the agent performs a random policy, and it is more likely to receive a penalty. It migrates outside of the flow domain or does not reach the destination. During the timesteps between 50 000 and 100 000, the agent has more chances to reach the destination and receives a positive reward. However, the agent may still fail the migration task due to additional explorations, which is a feature of the underlying governing reinforcement learning algorithm. For timesteps larger than 100 000,
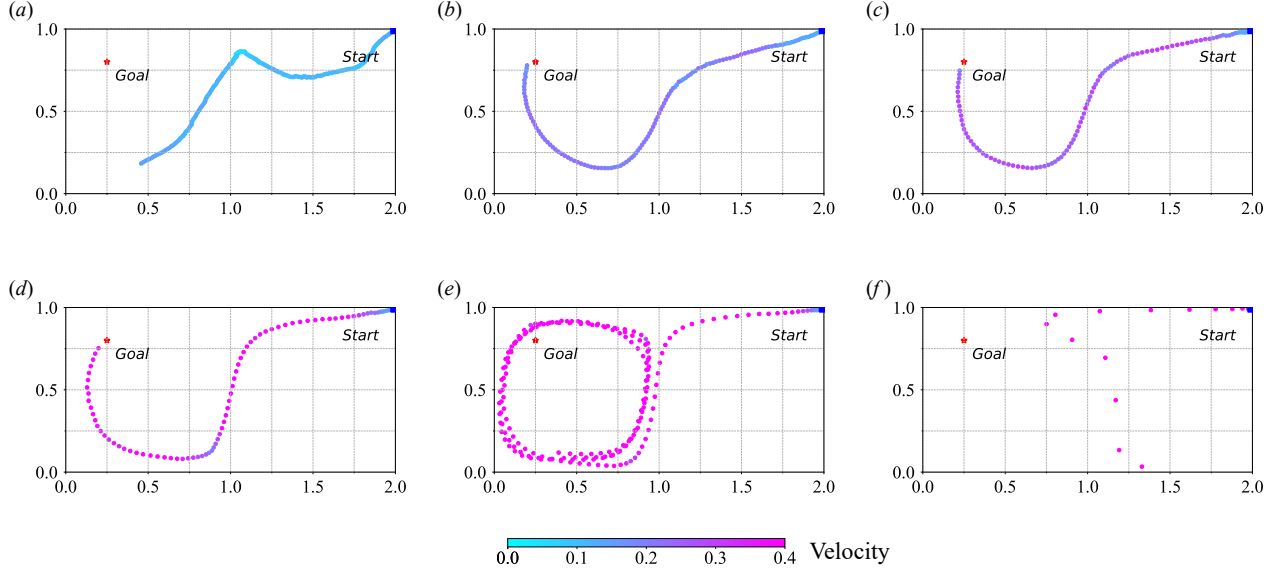
FIG. 4. Trajectories for the smart agent in the double-gyre flow with various control parameters of the flow field: (*a*) $A = 0.01$, (*b*) $A = 0.05$, (*c*) $A = 0.08$, (*d*) $A = 0.12$, (*e*) $A = 0.2$, and (*f*) $A = 1.0$. The trajectories are color-coded by instantaneous velocity magnitude.

the orange line gradually reaches a plateau with positive rewards, suggesting the training process converges. In Fig. 5, we also include trajectories for one failure model (with a reward around -5) and one successful model with moderate reward (around 5). For the failure model, we can see the agent only oscillates inside the right gyre, and its trajectory forms periodic orbits. For the successful model with moderate reward, the agent bypass yet misses the destination for the first time. It can then drift with the background flow and reach the destination for the second time.

The above results are obtained with the prescribed and fixed origin and destination positions for the agent. We further test the robustness of the energy-efficient strategies with respect to random positions of the origin and the destination. We choose the origin position in the right-half of the domain (i.e., $1 < x < 2$ and $0 < y < 1$) and the destination position in the left-half of the domain (i.e., $0 < x < 1$ and $0 < y < 1$). We performed 100 experiments with various origin and destination positions, and three typical optimized trajectories are shown in Figs. 6(a)-6(c). We can see that the smart agents always try to utilize the current of the carrier flow as much as possible. In addition, we calculate their accumulative energy consumption compared to that of the naive agents, and we plot the results in Figs. 6(d)-6(f). We can see that the smart agents still consume far less energy compared to the naive agents with the randomly chosen origin and destination positions.
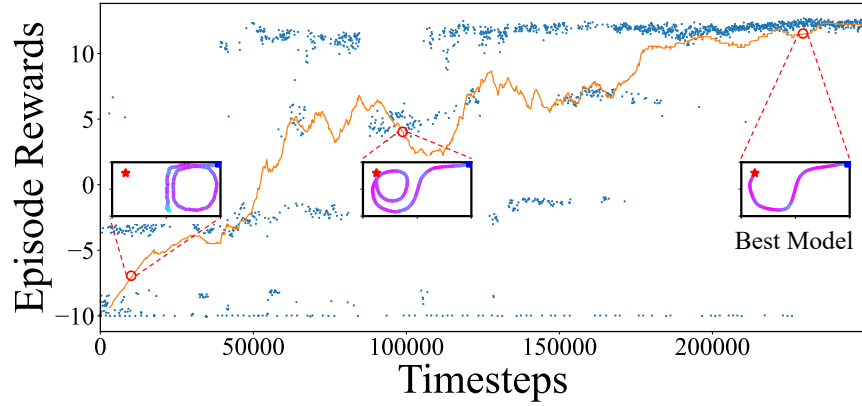
11

FIG. 5. The episode rewards as a function of timesteps during the training process. The discrete blue dots represent rewards obtained by the agents at the different episodes, and the orange line represents a smooth average of the rewards during a short-time window of 100 episodes. The insets include trajectories for one failure model, one successful model with moderate reward, and the successful model with the highest reward.
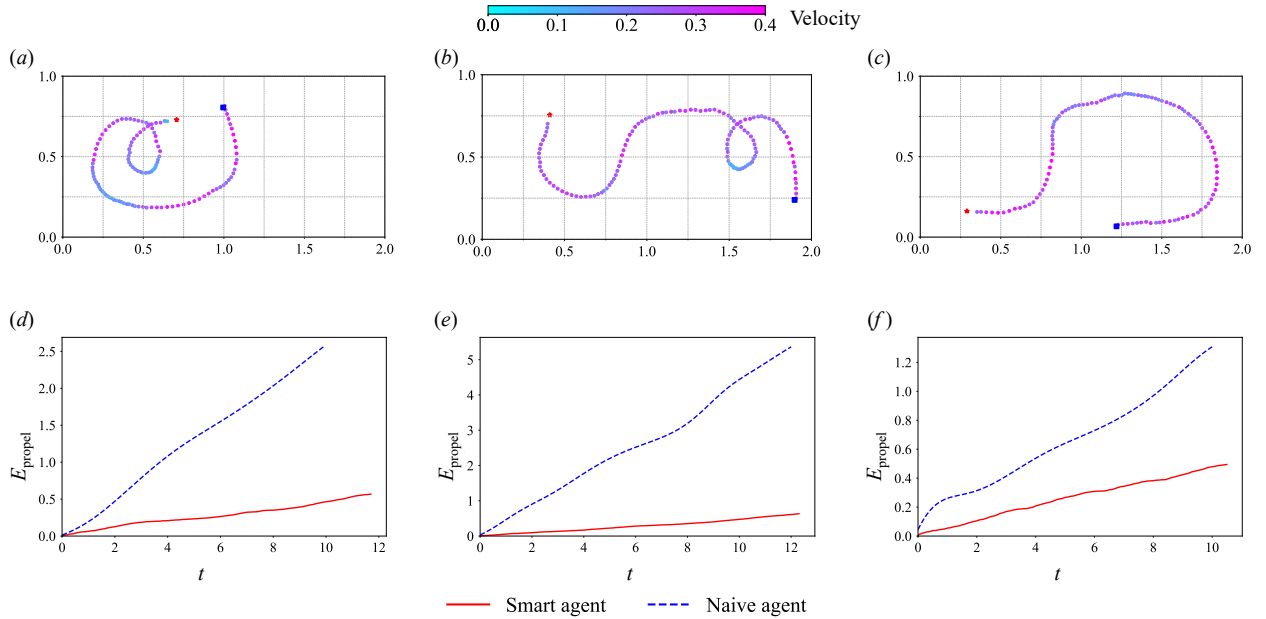


FIG. 6. (*a*)-(*c*) Trajectories for the smart agent in the double-gyre flow and (*d*)-(*f*) and the corresponding accumulative energy consumption compared to that of the naive agent.

## IV. MIGRATION IN THE TURBULENT RAYLEIGH-BÉNARD CONVECTION

### A. Numerical simulation of the turbulent Rayleigh-Bénard convection

We consider an incompressible thermal flow in the Oberbeck-Boussinesq approximation. The temperature is treated as an active scalar, and its influence on the velocity field is realized through the buoyancy term. The governing equations read as[22]

$$\nabla \cdot \mathbf{u} = 0 \tag{13a}$$

$$\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} = -\frac{1}{\rho_0}\nabla P + \nu \nabla^2 \mathbf{u} + g\beta_T(T - T_0)\hat{\mathbf{y}} \tag{13b}$$

$$\frac{\partial T}{\partial t} + \mathbf{u} \cdot \nabla T = \alpha_T \nabla^2 T \tag{13c}$$

where $\mathbf{u} = (u, v)$, $P$ and $T$ are velocity, pressure, and temperature of the fluid, respectively. $\rho_0$ and $T_0$ are the reference density and temperature, respectively. $\hat{\mathbf{y}}$ is the unit vector parallel to the gravity. $g$ is the gravity acceleration value. $\nu$, $\beta_T$, $\nu$, and $\alpha_T$ are the kinetic viscosity, thermal expansion coefficient, kinematic viscosity, and thermal diffusivity of the fluid, respectively. We adopt the lattice Boltzmann (LB) method as the numerical tool to solve the above equations. The advantages of the LB method include easy implementation and parallelization.[41–43] More numerical details on the LB method and validation of the in-house code can be found in our previous work.[44,45] In simulation, the top and bottom walls of the convection cell are kept at constant cold temperature $T_{\text{cold}}$ and hot temperature $T_{\text{hot}}$, respectively; while the other two vertical walls are adiabatic. All four walls impose no-slip velocity boundary conditions. The dimension of the cell is $L \times H$, and we set $L = 2H$ in this work. Simulation results are provided for the Prandtl number of $Pr = \nu/\alpha = 0.71$ and the Rayleigh number of $Ra = g\beta_T(T_{\text{hot}} - T_{\text{cold}})H^3/(\nu\alpha) = 10^8$. The $Pr$ corresponds to the thermal properties of air, while the $Ra$ is far less than that in the atmosphere due to the limitation of computational resources to simulate ultra-high Ra convection. Nevertheless, we can observe the large-scale coherent structure consisting of two primary rolls horizontally stacked in the simulation domain.[46–48]

### B. Training results and discussion

In training, we restrict the maximum propelling velocity of the agent along either $x$- or $y$-direction to be less than 0.02, which leads to $\|\mathbf{u}_{\text{propel}}\|^2 \leq 0.0008$. We show the instantaneous
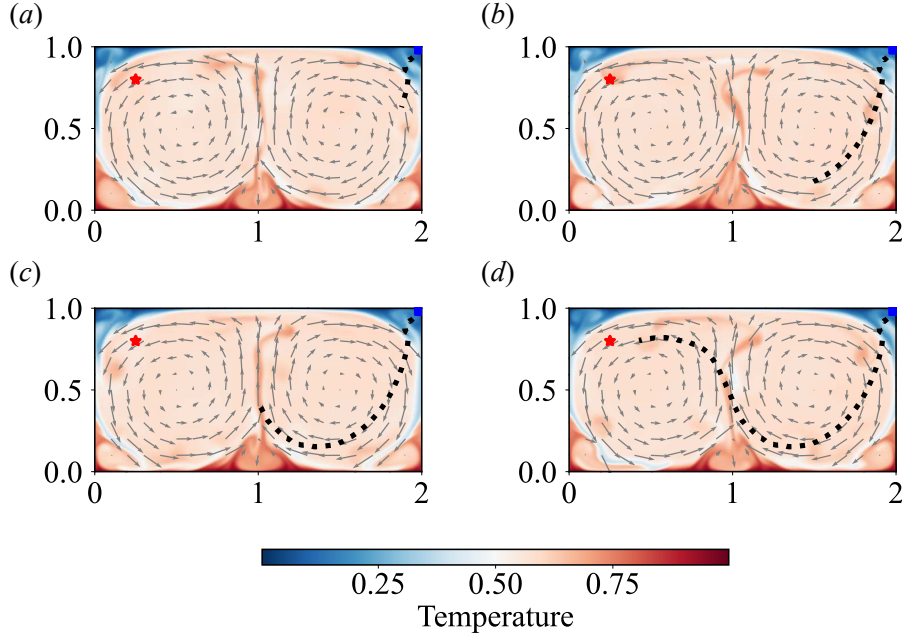
FIG. 7. Trajectory (black dotted line) of the smart agent in the two-dimensional turbulent Rayleigh-Bénard convection at (a) $t = 15$, (b) $t = 30$, (c) $t = 45$, and (d) $t = 65$. The contour shows the temperature field, and the vectors denote the velocity field of the convection. (Multimedia view)

trajectories for the smart agent in the turbulent Rayleigh-Bénard convection in Fig. 7 (Multimedia view). We released the agent at (2.0, 1.0), namely, the top-right corner in the domain (marked by the blue square in the plot). The agent's goal is to reach the destination position of (0.25, 0.8), marked by the red star in the plot with minimal energy consumption. Initially, the smart agent will move downward to utilize the vertical currents in the top-right region [see Fig. 7(a)]. When the agent reaches the bottom-right corner of the right roll, it will move leftward to utilize the horizontal currents [see Fig. 7(b)]. After that, it will drift into the bottom-right corner of the left roll [see Fig. 7(c)] and then follow the vertical upward current to achieve thermal soaring. When the agent reaches the same altitude as the destination, it will move leftward to utilize the horizontal current [see Fig. 7(d)]. The smart agent tries to follow the background currents as much as possible, which is similar to that in the double-gyre flow. In addition, the correlation coefficient between the orientation of the propelling velocity vector (i.e., $\theta_{propel}$) and the orientation of the fluid velocity vector (i.e., $\theta_{fluid}$) is 0.59, suggesting the statistical relevance between them.

In the turbulent Rayleigh-Bénard convection, we also compare the smart agent with the naive agent that moves straight from the origin to the destination. In Fig. 8, we show naive and smart
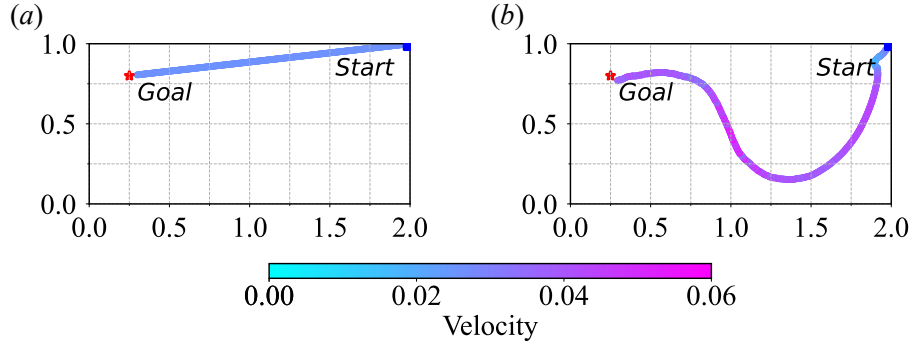
FIG. 8. Comparison of trajectories in the Rayleigh-Bénard convection for (*a*) a naive agent moving straightly and (*b*) a smart agent utilizing the flow structure to save energy. The trajectories are color-coded by the instantaneous velocity magnitude.

agents' trajectories, which are color-coded by the instantaneous velocity magnitude. We set the naive agent spending the same time as the smart agent to migrate from the origin to the destination. The naive agent migrates slower than the smart agent, because it travels a shorter distance than the smart agent. The trajectories differences for the smart agent in the double-gyre flow (see Fig 1) and the turbulent Rayleigh-Bénard convection (see Fig 7) are mainly due to the rotational direction of the primary vortex in the flows, while the underlying principles to utilize the flow structure to save energy remain the same for the smart agent.

We then plot the time series of accumulative energy consumed by the naive and smart agents in Fig. 9(a). We can see that both agents consume almost the same amount of energy during the initial period (i.e., around $t < 10$). The reason is that the background flow velocity is relatively small [i.e., around $O(10^{-5})$ near the origin position], and the migration of both agents heavily relies on their own propelling energy. After that (i.e., around $t > 10$), the smart agent moves downward to utilize the flow currents, while the naive agent continues to move straight toward the destination. Migrating reversely against the flow currents and crossing the edge of the roll will require substantial energy consumption, as evident from the significant increase in energy consumption at $t > 10$ [see Fig. 9(a)]. The results suggest that in turbulent flows with strong background flow velocity fluctuations, the smart agent can still save energy consumption while migrating to the destination. We also compare the accumulative total kinetic energy of the agents [see Fig. 9(b)]. After reaching the destination, the total kinetic energy of the smart agent is triple as that of the naive agent, while the smart agent consumed almost one-third of the energy.
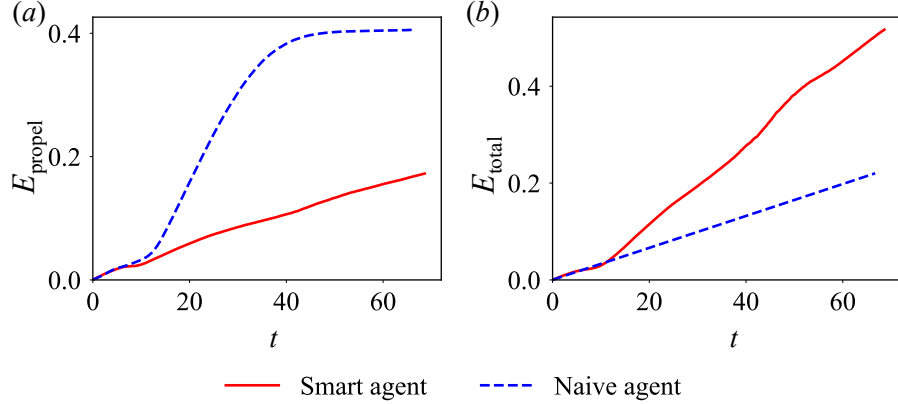
15

FIG. 9. Comparison of the accumulative (*a*) energy consumption $E_{\mathrm{propel}}$ and (*b*) total kinetic energy $E_{\mathrm{total}}$ of the smart agent and the naive agent in the Rayleigh-Bénard convection.

We show the episode rewards as a function of time steps during the training process in Fig. 10. The discrete blue dots represent rewards obtained by the agents at the different episodes. A reward value around -10 indicates that the agent migrates outside of the flow domain received a penalty; a reward value between -5 and 0 indicates that the agent neither migrates outside of the flow domain nor reaches the destination but is trapped in the right roll and oscillates inside the right roll; a reward value around 10 indicates that the agent can successfully reach the goal with minimal energy cost and migration time. We can see from Fig. 10 that the training processes for the agent in the turbulent flows can hardly converge in contrast to that in the simple periodic double-gyre flow. The main reason is that the flow field in the turbulent Rayleigh-Bénard convection exhibits more substantial fluctuations, and the migrating agent controlled by the reinforcement learning algorithm is more likely to explore different migration strategies in each episode. We determine the policy $\pi^*(\theta)$ that associated with the highest reward as the optimal policy, and we then train the agent with $\pi^*(\theta)$ to obtain the trajectory corresponding to the "best model". Due to the fluctuations of the flow field in the turbulent Rayleigh-Bénard convection, the trajectories generated by $\pi^*(\theta)$ may be slightly different in different runs. In Fig. 10, we include trajectories for two failure models (with a reward around -5 and -10, respectively) and one successful model with the highest reward (around 10). For the failure models, we can see the agent either oscillates in the right roll and does not drift into the left roll (around -5), or migrates outside of the domain (around -10).

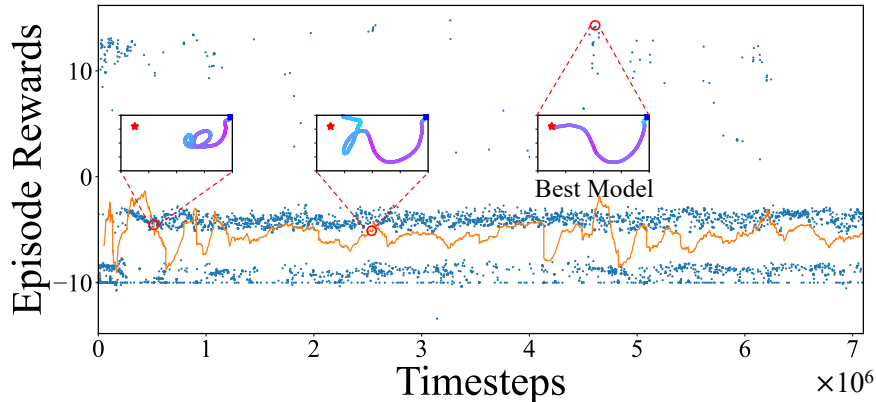The above results are obtained with observation variables including flow velocity, the agent's

16

FIG. 10. The episode rewards as a function of timesteps during the training process. The discrete blue dots represent rewards obtained by the agents at the different episodes, and the orange line represents a smooth average of the rewards during a short-time window of 100 episodes. The insets include trajectories for two failure models and the successful model with the highest reward.

spatial coordinates, and current time; however, in the turbulent Rayleigh-Bénard convection, the temperature acts as an active scalar that influences the velocity; thus, in Fig. 11 we further show the episode rewards in which the fluid temperature is also considered as additional observation variables during the training process. We can see that including temperature as sensorimotor, indeed, improves the average episode rewards during the training process, as evident there are more chances for the agent to obtain a high episode reward. On the other hand, we checked the optimized trajectories and the associated energy consumptions of the agents, and we found that including temperature as sensorimotor only slightly changes the optimized policy.

## V. CONCLUSIONS

In this work, we performed training of the self-propelling agent to migrate in a flow environment with the reinforcement learning algorithm. The smart agent that migrates in both the two-dimensional periodical double-gyre flow and two-dimensional turbulent Rayleigh-Bénard convection can learn to migrate from one position to another while utilizing the background flow currents as much as possible. We calculated the energy consumption for a naive agent that moves straight from the origin to the destination. The results show that the smart agent consumed less than one-fifth (or one-third) of energy compared to the naive agent in the periodical double-gyre flow (or the
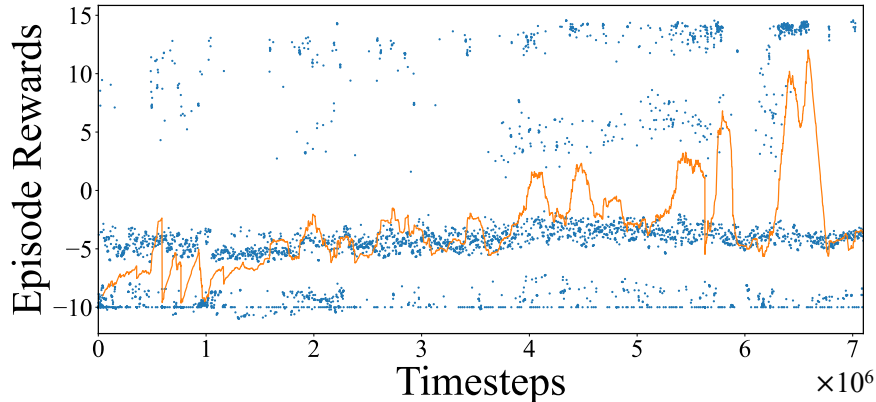
FIG. 11. The episode rewards as a function of timesteps during the training process. Different from the results presented in Fig. 10, here, we also include temperature as sensorimotor during the training process.

turbulent Rayleigh-Bénard convection). In addition, we found that compared to the double-gyre flow, the flow field in the turbulent Rayleigh-Bénard convection exhibits more substantial fluctuations, and the training agent is more likely to explore different migration strategies. Despite that the training process is challenging to converge for the smart agent in a turbulent environment, we can identify an energy-efficient trajectory that corresponds to the strategy with the highest reward received by the agent. We also found that including temperature as sensorimotor improves the average episode rewards during the training process in the turbulent Rayleigh-Bénard convection, while the optimized trajectories and the associated energy consumption of the agents almost remain unchanged. As pointed out by Laurent *et al.*,[40] there are opportunities to harness the energy of turbulence, particularly for person-less transport and small reconnaissance aircraft. Thus, similar processes could very well be optimized in other migration involving a self-propelling agent, such as UAVs flying in a turbulent convective environment.

## ACKNOWLEDGMENTS

## AUTHOR DECLARATIONS

### Conflict of Interest

The authors have no conflicts to disclose.

## DATA AVAILABILITY

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

## APPENDIX: PROPULSION VELOCITY VECTOR OF THE SMART AGENT IN THE DOUBLE-GYRE FLOW

We show the propulsion velocity vector of the smart agent in different locations in Fig. 12, which further clarify the dynamics of the smart agent in the double-gyre flow. The smart agent utilizes background currents as much as possible, as evident that the angle between the propelling velocity vector (the black vector) and the fluid velocity vector (color-coded vectors) is generally less than $90°$ at the same location.

## REFERENCES

[1]C. D. Cone, "Thermal soaring of birds," American Scientist **50**, 180–209 (1962).

[2]F. Ludlam and R. Scorer, "Reviews of modern meteorology 10 convection in the atmosphere," Quarterly Journal of the Royal Meteorological Society **79**, 317–341 (1953).

[3]R. Bencatel, J. T. de Sousa, and A. Girard, "Atmospheric flow field models applicable for aircraft endurance extension," Progress in Aerospace Sciences **61**, 1–25 (2013).

[4]H. J. Williams, E. Shepard, M. D. Holton, P. Alarcón, R. Wilson, and S. Lambertucci, "Physical limits of flight performance in the heaviest soaring bird," Proceedings of the National Academy of Sciences **117**, 17884–17890 (2020).

[5]P. B. MacCready, "Optimum airspeed selector," Soaring (January–February) **10**, 10 (1958).

[6]Z. Akos, M. Nagy, and T. Vicsek, "Comparing bird and human soaring strategies," Proceedings of the National Academy of Sciences **105**, 4139–4143 (2008).
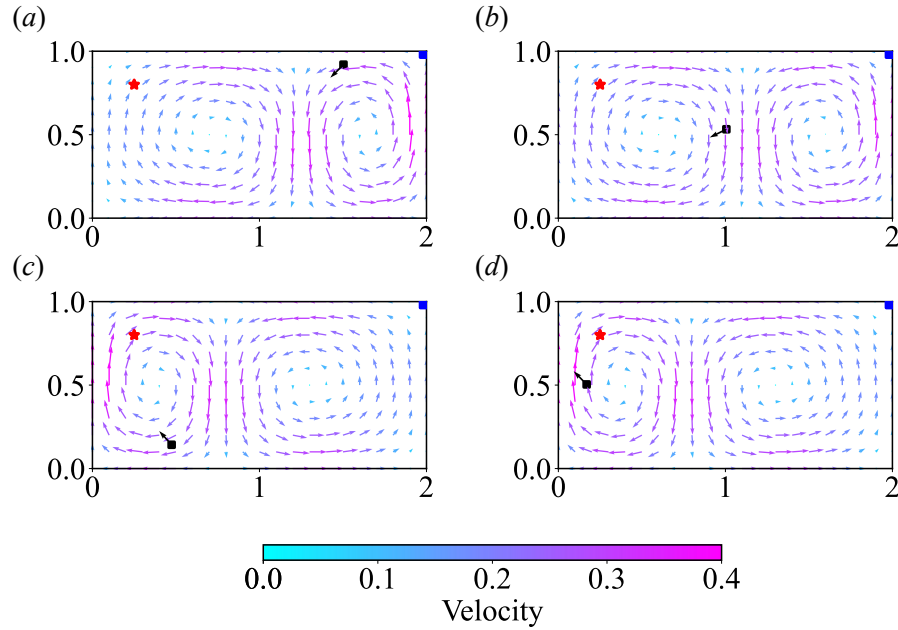
FIG. 12. Propulsion velocity vector (black vector) of the smart agent in the double-gyre flow at (*a*) $t = 2.2$, (*b*) $t = 4.2$, (*c*) $t = 6.7$, and (*d*) $t = 8.2$. The color-coded vectors denote the velocity field of the double-gyre flow.

[7]M. J. Allen and V. Lin, "Guidance and control of an autonomous soaring uav," Tech. Rep. (2007).

[8]J. Wharington and I. Herszberg, "Control of a high endurance unmanned air vehicle," in *Proceedings of the 21st ICAS Congress*, Vol. 1234567890 (1998).

[9]Z. Ákos, M. Nagy, S. Leven, and T. Vicsek, "Thermal soaring flight of birds and unmanned aerial vehicles," Bioinspiration & Biomimetics **5**, 045003 (2010).

[10]G. Reddy, A. Celani, T. J. Sejnowski, and M. Vergassola, "Learning to soar in turbulent environments," Proceedings of the National Academy of Sciences **113**, E4877–E4884 (2016).

[11]G. Reddy, J. Wong-Ng, A. Celani, T. J. Sejnowski, and M. Vergassola, "Glider soaring via reinforcement learning in the field," Nature **562**, 236–239 (2018).

[12]T. Dbouk and D. Drikakis, "Quadcopter drones swarm aeroacoustics," Physics of Fluids **33**, 057112 (2021).

[13]S. Colabrese, K. Gustavsson, A. Celani, and L. Biferale, "Flow navigation by smart microswimmers via reinforcement learning," Physical Review Letters **118**, 158004 (2017).

[14]K. Gustavsson, L. Biferale, A. Celani, and S. Colabrese, "Finding efficient swimming strategies in a three-dimensional chaotic flow by reinforcement learning," The European Physical Journal E **40**, 1–6 (2017).

20

[15] S. Colabrese, K. Gustavsson, A. Celani, and L. Biferale, "Smart inertial particles," Physical Review Fluids **3**, 084301 (2018).

[16] L. Biferale, F. Bonaccorso, M. Buzzicotti, P. Clark Di Leoni, and K. Gustavsson, "Zermelo's problem: Optimal point-to-point navigation in 2D turbulent flows using reinforcement learning," Chaos: An Interdisciplinary Journal of Nonlinear Science **29**, 103138 (2019).

[17] J. K. Alageshan, A. K. Verma, J. Bec, and R. Pandit, "Machine learning strategies for path-planning microswimmers in turbulent flows," Physical Review E **101**, 043110 (2020).

[18] S. Muiños-Landin, A. Fischer, V. Holubec, and F. Cichos, "Reinforcement learning with artificial microswimmers," Science Robotics **6**, eabd9285 (2021).

[19] J. Qiu, W. Huang, C. Xu, and L. Zhao, "Swimming strategy of settling elongated micro-swimmers by reinforcement learning," SCIENCE CHINA Physics, Mechanics & Astronomy **63**, 1–9 (2020).

[20] J. Qiu, N. Mousavi, K. Gustavsson, C. Xu, B. Mehlig, and L. Zhao, "Navigation of micro-swimmers in steady flow: The importance of symmetries," Journal of Fluid Mechanics **932**, A10 (2022).

[21] S. C. Shadden, F. Lekien, and J. E. Marsden, "Definition and properties of Lagrangian coherent structures from finite-time Lyapunov exponents in two-dimensional aperiodic flows," Physica D: Nonlinear Phenomena **212**, 271–304 (2005).

[22] K.-Q. Xia, "Current trends and future directions in turbulent thermal convection," Theoretical and Applied Mechanics Letters **3**, 052001 (2013).

[23] F. M. Callier and C. A. Desoer, *Linear system theory* (Springer, 2012).

[24] L. S. Pontryagin, *Mathematical theory of optimal processes* (CRC press, 1987).

[25] P. Mehta, M. Bukov, C.-H. Wang, A. G. Day, C. Richardson, C. K. Fisher, and D. J. Schwab, "A high-bias, low-variance introduction to Machine Learning for physicists," Physics Reports **810**, 1–124 (2019).

[26] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, "Machine learning and the physical sciences," Reviews of Modern Physics **91**, 045002 (2019).

[27] S. L. Brunton, B. R. Noack, and P. Koumoutsakos, "Machine learning for fluid mechanics," Annual Review of Fluid Mechanics **52**, 477–508 (2020).

[28] 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).

[29] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction* (MIT press, 2018).

[30] F. Cichos, K. Gustavsson, B. Mehlig, and G. Volpe, "Machine learning for active matter," Nature Machine Intelligence **2**, 94–103 (2020).

[31] J. N. Tsitsiklis and B. Van Roy, "An analysis of temporal-difference learning with function approximation," IEEE Transactions on Automatic Control **42**, 674–690 (1997).

[32] C. Szepesvári, "Algorithms for reinforcement learning," Synthesis Lectures on Artificial Intelligence and Machine Learning **4**, 1–103 (2010).

[33] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International Conference on Machine Learning* (PMLR, 2018) pp. 1861–1870.

[34] W. Zhang, T. Inanc, S. Ober-Blobaum, and J. E. Marsden, "Optimal trajectory generation for a glider in time-varying 2D ocean flows B-spline model," in *2008 IEEE International Conference on Robotics and Automation* (IEEE, 2008) pp. 1083–1088.

[35] C. Wang, Z. Xu, X. Zhang, and S. Wang, "Optimal reduced frequency for the power efficiency of a flat plate gliding with spanwise oscillations," Physics of Fluids **33**, 111908 (2021).

[36] L. Liu, M. Chen, J. Yu, Z. Zhang, and X. Wang, "Full-scale simulation of self-propulsion for a free-running submarine," Physics of Fluids **33**, 047103 (2021).

[37] W. Wang, H. Huang, and X.-Y. Lu, "Optimal chordwise stiffness distribution for self-propelled heaving flexible plates," Physics of Fluids **32**, 111905 (2020).

[38] H. Yu, X.-Y. Lu, and H. Huang, "Collective locomotion of two uncoordinated undulatory self-propelled foils," Physics of Fluids **33**, 011904 (2021).

[39] Y. Liu, C. Pan, and Y. Liu, "Propulsive performance and flow-field characteristics of a jellyfish-like ornithopter with asymmetric pitching motion," Physics of Fluids **32**, 071904 (2020).

[40] K. M. Laurent, B. Fogg, T. Ginsburg, C. Halverson, M. J. Lanzone, T. A. Miller, D. W. Winkler, and G. P. Bewley, "Turbulence explains the accelerations of an eagle in natural flight," Proceedings of the National Academy of Sciences **118**, e2102588118 (2021).

[41] A. Xu, W. Shyy, and T. Zhao, "Lattice Boltzmann modeling of transport phenomena in fuel cells and flow batteries," Acta Mechanica Sinica **33**, 555–574 (2017).

[42] S. Chen and G. D. Doolen, "Lattice Boltzmann method for fluid flows," Annual Review of Fluid Mechanics **30**, 329–364 (1998).

[43] C. K. Aidun and J. R. Clausen, "Lattice-Boltzmann method for complex flows," Annual Review of Fluid Mechanics **42**, 439–472 (2010).

[44] A. Xu, L. Shi, and T. Zhao, "Accelerated lattice Boltzmann simulation using GPU and OpenACC with data management," International Journal of Heat and Mass Transfer **109**, 577–588 (2017).

[45] A. Xu, L. Shi, and H.-D. Xi, "Lattice Boltzmann simulations of three-dimensional thermal convective flows at high Rayleigh number," International Journal of Heat and Mass Transfer **140**, 359–370 (2019).

[46] W.-F. Zhou and J. Chen, "Large-scale structures of turbulent Rayleigh–Bénard convection in a slim-box," Physics of Fluids **33**, 065103 (2021).

[47] A. Xu, X. Chen, F. Wang, and H.-D. Xi, "Correlation of internal flow structure with heat transfer efficiency in turbulent Rayleigh–Bénard convection," Physics of Fluids **32**, 105112 (2020).

[48] A. Xu, B.-R. Xu, L.-S. Jiang, and H.-D. Xi, "Production and transport of vorticity in two-dimensional Rayleigh-Bénard convection cell," Physics of Fluids **34**, 013609 (2022).