# Efficient collective swimming by harnessing vortices through deep reinforcement learning

Siddhartha Verma [1†], Guido Novati [1†], Petros Koumoutsakos[1*]

[1]*Computational Science and Engineering Laboratory, Clausiusstrasse 33, ETH Zürich, CH-8092, Switzerland*

[*]*Corresponding author e-mail: petros@ethz.ch*

**Classification**:

1. Physical Sciences/Applied Mathematics
2. Biological Sciences/Biophysics and Computational Biology

**Keywords**: Fish schooling; Deep reinforcement learning; Autonomous navigation

**Short Title**: Efficient propulsion through deep reinforcement learning

---

[†]Authors contributed equally to this work

## Abstract

Fish in schooling formations navigate complex flow-fields replete with mechanical energy in the vortex wakes of their companions. Their schooling behaviour has been associated with evolutionary advantages including collective energy savings. How fish harvest energy from their complex fluid environment and the underlying physical mechanisms governing energy-extraction during collective swimming, is still unknown. Here we show that fish can improve their sustained propulsive efficiency by actively following, and judiciously intercepting, vortices in the wake of other swimmers. This swimming strategy leads to collective energy-savings and is revealed through the first ever combination of deep reinforcement learning with high-fidelity flow simulations. We find that a 'smart-swimmer' can adapt its position and body deformation to synchronise with the momentum of the oncoming vortices, improving its average swimming-efficiency at no cost to the leader. The results show that fish may harvest energy deposited in vortices produced by their peers, and support the conjecture that swimming in formation is energetically advantageous. Moreover, this study demonstrates that deep reinforcement learning can produce navigation algorithms for complex flow-fields, with promising implications for energy savings in autonomous robotic swarms.

**Significance Statement.** *Fish schooling is one of the most intriguing instances of collective behavior and complex decision making in nature, yet its underlying physical mechanisms remain largely unknown. We combine state of the art flow simulations with reinforcement learning, to answer the longstanding question of whether schooling fish may reduce energy-expenditure by adapting their swimming motion to the flow created by their companions. We demonstrate that a 'smart' self-propelled swimmer can autonomously adapt its swimming behaviour to exploit energy deposited in the wake of other swimmers. The results support the thesis that fish may exploit unsteady flow-fields generated by collective locomotion to reap substantial energetic benefits and have promising implications for autonomous robotic swarms.*

There is a long-standing interest for understanding and exploiting the physical mechanisms employed by active swimmers in nature (nektons).[1–4] Fish schooling in particular, one of the most striking patterns of collective behaviour, has been the subject of intense investigations.[5–9] A key issue in understanding fish schooling behaviour, and its potential for engineering applications,[10] is the clarification of the role of the flow environment. Fish sense and navigate in complex flow-fields full of mechanical energy that is distributed across multiple scales by vortices generated by obstacles and other swimming organisms [11,12]. There is evidence that their swimming behaviour adapts to flow gradients (rheotaxis) and, in certain cases, it reflects energy-harvesting from such environments.[13–15] Hydrodynamic interactions have also been implicated in the fish schooling patterns that form when individual fish adapt their motion to that of their peers, while compensating for flow-induced displacements. Recent experimental studies have argued that fish may interact beneficially with each other,[9,16,17] but in ways that challenge[18] the earlier proposed mechanisms [5,6] governing fish-schooling. However, the role of hydrodynamics in fish schooling is not embraced universally[8,19,20] and there is limited quantitative information regarding the physical mechanisms that would explain such energetic benefits. Experimental[16,17] and computational[21] studies of collective swimming have been hampered by the presence of multiple deforming bodies and their interactions with the flow-field. Moreover, numerical simulations have demonstrated that a coherent swimming group cannot be sustained without exerting some form of control strategy on the swimmers.[22,23] Here, we employ deep reinforcement learning (deep RL[24]) to discover such strategies for two autonomous and self-propelled swimmers, and elucidate the physical mechanisms that enable efficient and sustained coordinated swimming.

During fish propulsion, body undulations and the sideways displacement of the caudal fin generate and inject a series of vortex rings in its wake.[25–27] When fish swim in formation, these vortices may assist the locomotion of fish that intercept them judiciously, which in turn can reduce the collective swimming effort. Such vortex-induced benefits have been observed in trout, which curtail muscle usage by capitalizing on energy injected in the flow by obstacles present in streams.[14,28] Here we examine two self-propelled swimmers in a leader/follower arrangement, and investigate the physical mechanisms that lead to energetically beneficial interactions by considering four distinct scenarios. Two of these involve smart-followers that can take autonomous decisions when interacting with a leader's wake, and are referred to as Interacting Swimmers ($IS$) (e.g., the follower in Fig. 1). Additionally, we consider two distinct Solitary Swimmers ($SS$) that swim in isolation in an unbounded domain. In the case of interacting swimmers, $IS_\eta$ denotes swimmers that learn the most efficient way of swimming in the leader's wake (without any positional constraints), and acquire a policy $\pi_\eta$ in the process. In turn, swimmer $IS_d$ attempts to
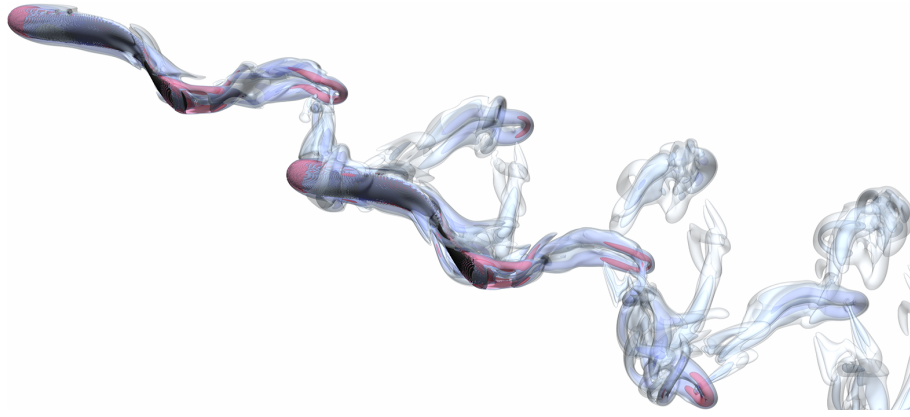
minimize lateral deviations from the leader's path, resulting in optimal policy $\pi_d$. These autonomous swimmers take decisions by virtue of deep RL, using visual cues from their environment (see Fig. 1d). The Solitary Swimmers $SS_\eta$ and $SS_d$ execute actions identical to $IS_\eta$ and $IS_d$, respectively, and serve as 'control' configurations to assess how the absence of a leader's wake impacts swimming-energetics.

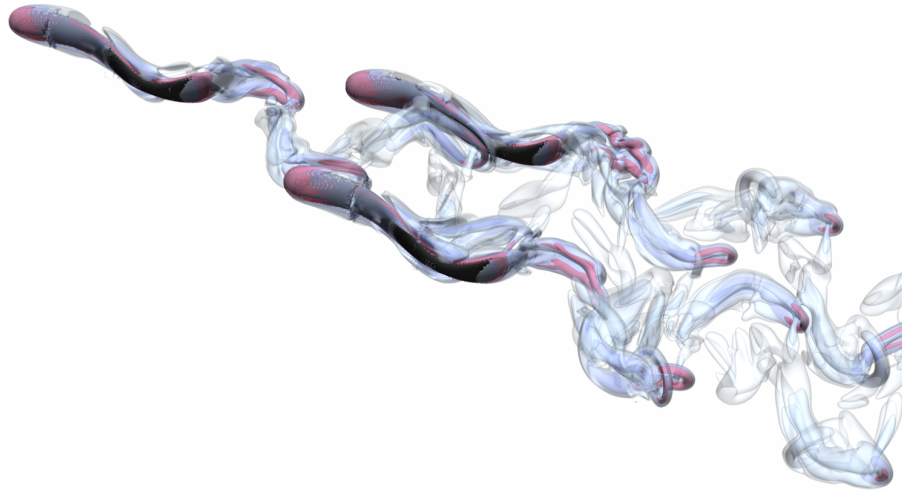## Reinforcement learning for autonomous swimmers

Reinforcement learning [29] has been introduced to identify navigation policies in several model systems of vortex dipoles, soaring birds and micro-swimmers.[30–32] Here, we expand on our earlier work[22,33] combining Reinforcement Learning with Direct Numerical Simulations of the Navies stokes equations for two self-propelled and autonomous swimmers. We first investigate two-dimensional swimmers in a tandem configuration and analyse their kinematics for the cases of $IS_\eta$ and $IS_d$ (Fig. 2). In both cases, the swimmer trails a leader representing an adult zebrafish of length $L$, swimming steadily at a velocity $U$ (Reynolds number $Re = UL/\nu \approx 5000$). We employ deep Reinforcement Learning (see Methods section for details), and after training we observe that $IS_d$ is able to maintain its position behind the leader quite effectively ($\Delta y \approx 0$, Fig. 2b), in accordance to its reward ($R_d = 1 - |\Delta y|/L$). Surprisingly, $IS_\eta$ with a reward function proportional to swimming-efficiency ($R_\eta = \eta$), also settles close to the center of the leader's wake (Fig. 2b and Supplementary Movie S2), although it receives no reward associated with its relative position. Both $IS_d$ and $IS_\eta$ maintain a distance of $\Delta x \approx 2.2L$ from their respective leaders (Figure 2a). $IS_\eta$ shows a greater proclivity to maintain this separation and intercepts the periodically shed wake-vortices just after they have been fully formed and detach from the leader's tail. In addition to $\Delta x = 2.2L$, there is an additional point of stability at $\Delta x = 1.5$ (Fig. 2c). The difference $0.7L$ matches the distance between vortices in the wake of the leader. In both positions the lateral motion of the follower's head is synchronized with the flow-velocity in the leader's wake, thus inducing minimal disturbance on the oncoming flow-field. We note that a similar synchronization has been observed when trout minimize muscle usage by interacting with vortex-columns in a cylinder's wake.[14] $IS_\eta$ undergoes relatively minor body-deformation while manoeuvring (Figure 2d), whereas $IS_d$ executes aggressive turns involving large body-curvature. Trout interacting with cylinder-wakes exhibit increased body-curvature,[28] which is contrary to the behaviour displayed by $IS_\eta$. The difference may be ascribed to the widely-spaced vortex columns generated by large-diameter cylinders used in the experimental study. Weaving in and out of comparatively smaller vortices generated by like-sized fish encountered in a school (Fig. 1c) would entail excessive energy consumption. We note that maintaining $\Delta y = 0$ requires significant effort by $IS_d$ (Supplementary Fig. S2d) since its reward ($R_d$) is insensitive to energy expenditure. A previous study[33] suggested that minimizing lateral displacement led to enhanced swimming-efficiency (compared to the leader), albeit with noticeable deviation from $\Delta y = 0$. In the current study, recurrent neural networks augmented with 'Long Short-Term Memory' cells (Supplementary Fig. S3) help to encode time-dependencies in the value function, and enable far more robust smart-swimmers. Thus, stringent attempts by $IS_d$ to correct for oscillations about $\Delta y = 0$ (Fig. 2b) give rise to increased costs (Supplementary Fig. S2).

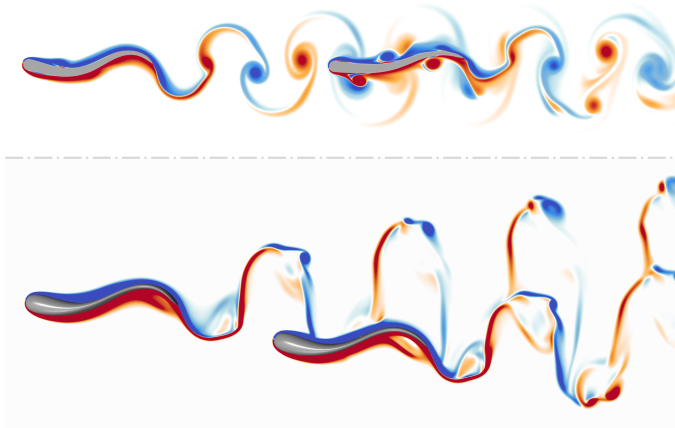## Intercepting vortices for efficient swimming

To determine the impact of wake-induced interactions on swimming-performance, we compare energetics data for $IS_\eta$ and $SS_\eta$ (Fig. 3). The swimming-efficiency of $IS_\eta$ is significantly higher than that of $SS_\eta$ (Fig. 3a), whereas the Cost of Transport (CoT), which represents energy spent for traversing a unit distance, is lower (Fig. 3b). Over a duration of 10 tail-beat periods (from $t = 20$ to $t = 30$, Supplementary Fig. S2) $IS_\eta$ experiences a 11% increase in average speed compared to $SS_\eta$, a 32% increase in average swimming-efficiency, and a 36% decrease in CoT. The benefit for $IS_\eta$ results from both a 29% reduction in effort required for deforming its body against flow-induced forces ($P_{Def}$), and a 53% increase in average thrust-power ($P_{Thrust}$). Performance-differences between $IS_\eta$ and $SS_\eta$ exist solely due to the presence/absence of a preceding wake, since both swimmers undergo identical body-undulations throughout the simulations. Comparing the swimming-efficiency and power values of four distinct swimmers (Supplementary Fig. S2 and Supplementary Table 1), we confirm that $IS_\eta$ and $SS_\eta$
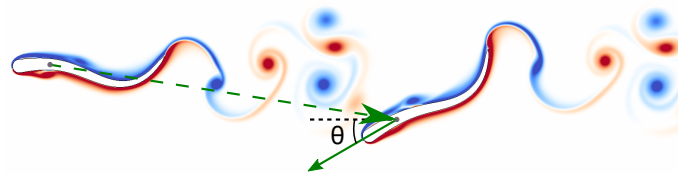
(a)



(b)



(c)



(d)

Figure 1: **Fluid-mediated interaction, and visual cues for a smart swimmer.** (a) 3D simulation of two nonautonomous swimmers, in which the leader swims steadily, and the follower maintains a specified relative position such that it interacts favourably with the leader's wake. (b) 3D simulation of three nonautonomous swimmers, where the two followers maintain specified relative positions that are beneficial. The flow-structures have been visualized using iso-surfaces of the Q-criterion. The vortex rings shed by each swimmer spread out in a diverging V-shaped pattern due to their self-induced velocity. An animation of the 3D simulation is provided in Supplementary Movie S1. (c) Comparison of vorticity field in the wake of 2D (top panel) and 3D (bottom panel) swimmers (red: positive, blue: negative). Every half a tail-beat period, the smart-follower in 2D simulations ($IS_\eta$) autonomously selects the most appropriate action encoded in policy $\pi_\eta$ learned during training-simulations, which allows it to maximize long-term swimming-efficiency (Supplementary Movie S2). The smart-follower is capable of adapting to deviations in the leader's trajectory (Supplementary Movie S3), as these situations are encountered when performing random actions during training. (d) The smart-swimmer relies on a pre-defined set of variables to identify its 'observed-state', some of which are depicted in this figure. Additional observed-state parameters are described in the Methods section.
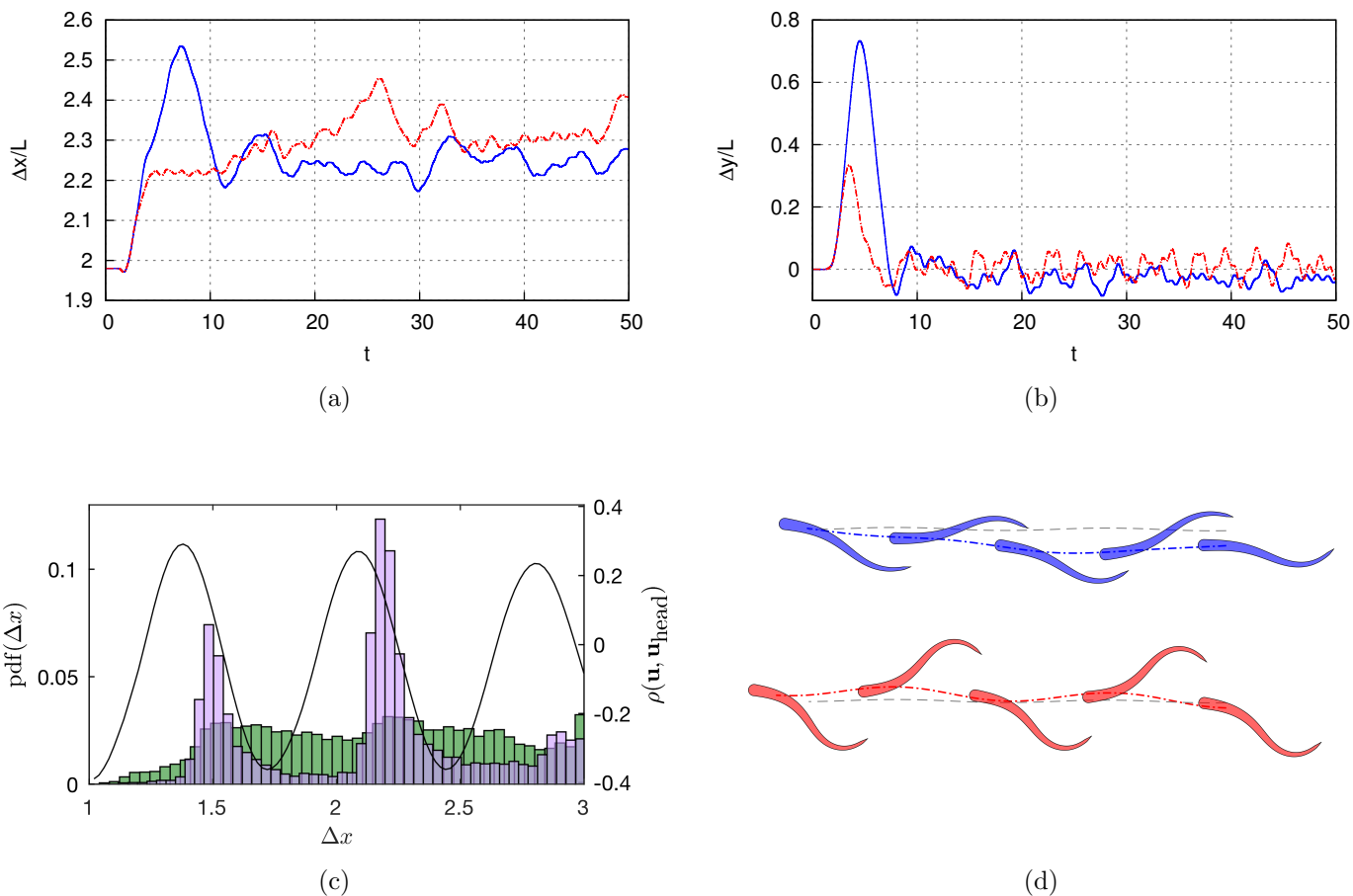
(a)

(b)





(c)

(d)

Figure 2: **Relative position, correlation with the flow-field, and severity of body-deformation.** (a) Relative horizontal displacement of the smart followers with respect to the leader, over a duration of 50 tail-beat periods starting from rest (solid blue line - $IS_\eta$, dash-dot red line - $IS_d$) (b) Lateral displacement of the smart followers. (c) Histogram showing the probability density function (pdf - left vertical axis) of swimmer $IS_\eta$'s preferred center-of-mass location during training. In the early stages of training (first 10000 transitions - green bars), the swimmer does not show a strong preference for maintaining any particular separation distance. Towards the end of training (last 10000 transitions - lilac bars), the swimmer displays a strong preference for maintaining a separation-distance of either $\Delta x = 1.5L$ or $2.2L$. The solid black line in the figure depicts correlation-coefficient, with peaks in the black curve signifying locations where the smart-follower's head-movement would be synchronized with the flow-velocity in an undisturbed wake (please see Supplementary Information for relevant details). (d) Comparison of body-deformation for swimmers $IS_\eta$ (top) and $IS_d$ (bottom), from $t = 27$ to $t = 29$. Their respective trajectories are shown with the dash-dot lines, whereas the dashed gray line represents the trajectory of the leader (not shown). A quantitative comparison of body-curvature for the two swimmers may be found in Supplementary Fig. S1.

are considerably more energetically efficient than either $IS_d$ or $SS_d$, thus verifying the hydrodynamic benefits of coordinated swimming.

The efficient swimming of $IS_\eta$ (e.g., point $\eta_{max}(A)$ in Fig. 3a) is attributed to the synchronized motion of its head with the lateral flow-velocity generated by the wake-vortices of the leader (see panel 'v' in Supplementary Movie S2). This mechanism is evidenced by the correlation-curve shown in Fig. 2c, and by the co-alignment of velocity vectors close to the head in Figs. 4a and 4b. As shown in Supplementary Movie S4, $IS_\eta$ intercepts the oncoming vortices in a slightly skewed manner, splitting each vortex into a stronger ($W_{1U}$, Fig. 4a) and a weaker fragment ($W_{1L}$). The vortices interact with the swimmer's own boundary layer to generate 'lifted-vortices' ($L_1$), which in turn generate secondary-vorticity ($S_1$) close to the body. Meanwhile, the wake- and lifted-vortices created during the previous half-period, $W_{2U}$, $W_{2L}$, and $L_2$, have travelled downstream along the body. This sequence of events alternates periodically between the upper (right-lateral) and lower (left-lateral) surfaces, as seen in Supplementary Movie S4. Interactions of $IS_\eta$ with the flow-field at points $\eta_{min}(D)$ and $(E)$ in Fig. 3a are analyzed separately in Supplementary Figs. S4 and S5.

We observe that the swimmer's upper surface is covered in a layer of negative vorticity (and vice versa for the lower surface) (Fig. 4a, top panel) owing to the no-slip boundary condition. The wake- or the lifted-vortices weaken this distribution by generating vorticity of opposite sign (e.g., secondary-vorticity visible in narrow regions between the fish-surface and vortices $L_1$, $W_{1L}$, $L_2$, and $L_3$), and create high-speed areas visible as bright spots in Fig. 4a (lower panel). The resulting low-pressure region exerts a suction-force on the surface of the swimmer (Fig. 4b, upper panel), which assists body-undulations when the force-vectors coincide with the deformation-velocity (Fig. 4b lower panel), or increases the effort required when they are counter-aligned. The detailed impact of these interactions is demonstrated in Figs. 4c to 4f. On the lower surface, $W_{1L}$ generates a suction-force oriented in the same direction as the deformation-velocity ($0 < s < 0.2L$ in Fig. 4b), resulting in negative $P_{Def}$ (Fig. 4e) and favourable $P_{Thrust}$ (Fig. 4f). On the upper surface, the lifted-vortex $L_1$ increases the effort required for deforming the body (positive peak in Fig. 4c at $s = 0.2L$), but is beneficial in terms of producing large positive thrust-power (Fig. 4d). Moreover, as $L_1$ progresses along the body, it results in a prominent reduction in $P_{Def}$ over the next half-period, similar to the negative peak produced by the lifted-vortex $L_2$ ($s = 0.55L$ in Fig. 4e). The average $P_{Def}$ on both the upper and lower surfaces is predominantly negative (i.e., beneficial), in contrast to the minimum swimming-efficiency instance $\eta_{min}(D)$, where a mostly positive $P_{Def}$ distribution signifies substantial effort required for deforming the body (Supplementary Fig. S4). We observe noticeable drag on the upper surface close to $s = 0$ (Fig. 4b top panel and Fig. 4d), attributed to high-pressure region forming in front of the swimmer's head. Forces induced by $W_{1L}$ are both beneficial and detrimental in terms of generating thrust-power ($0 < s < 0.2L$ in Fig. 4f), whereas forces induced by $L_2$ primarily increase drag but assist in body-deformation (Fig. 4e). The tail-section ($s = 0.8L$ to $1L$) does not contribute noticeably to either thrust- or deformation-power at the instant of maximum swimming-efficiency.

## Energy-saving mechanisms in coordinated swimming

The most discernible behaviour of $IS_\eta$ is the synchronization of its head-movement with the wake-flow. However, the most prominent reduction in deformation-power occurs near the midsection of the body ($0.4 \leq s \leq 0.7$ in Figs. 4c and 4e). This indicates that the technique devised by $IS_\eta$ is markedly different from energy-conserving mechanisms implied in previous theoretical[6,34] and computational[21] work, namely, drag-reduction attributed to reduced relative-velocity in the flow, and thrust-increase owing to the 'chanelling effect'. In fact, the predominant energetics-gain (i.e., negative $P_{Def}$) occurs in areas of high relative-velocity, for instance near the high-velocity spot generated by vortex $L_2$ (Fig. 4). This dependence of swimming-efficiency on a complex interplay between wake-vortices and body-deformation aligns closely with experimental findings.[14,28]

We remark that the majority of the results presented here were obtained with a steadily-swimming leader. However, with no additional training, $IS_\eta$ is able to extract an energetic-benefit even when exposed to an erratic leader (as seen in Supplementary Movie S3), where it deliberately chooses to interact with the unsteady wake. Moreover, given the head-synchronization tendency of the 2D smart-swimmer, we identify suitable locations behind a 3D leader where the flow velocity would match a follower's head motion (Supplementary Fig. S6). A feedback controller
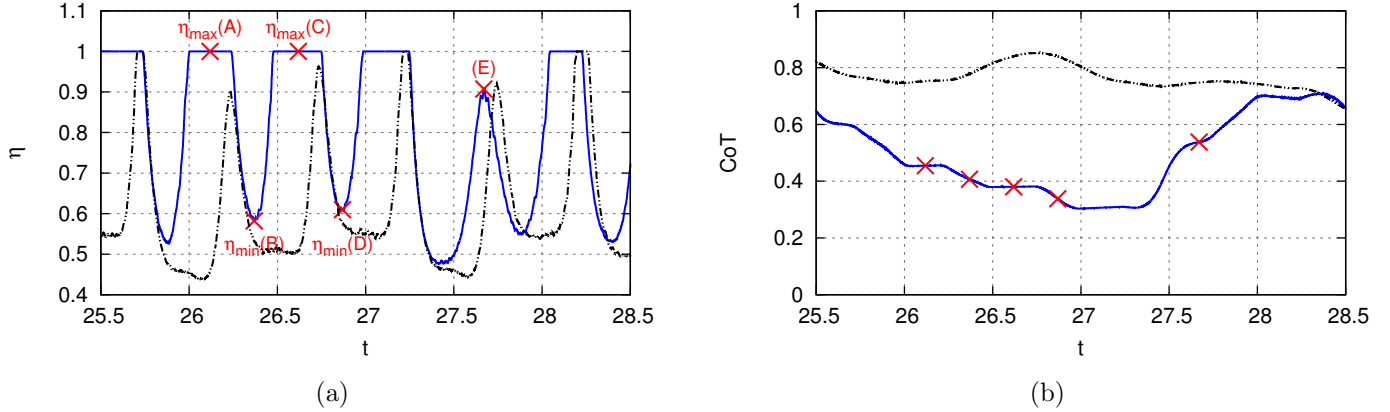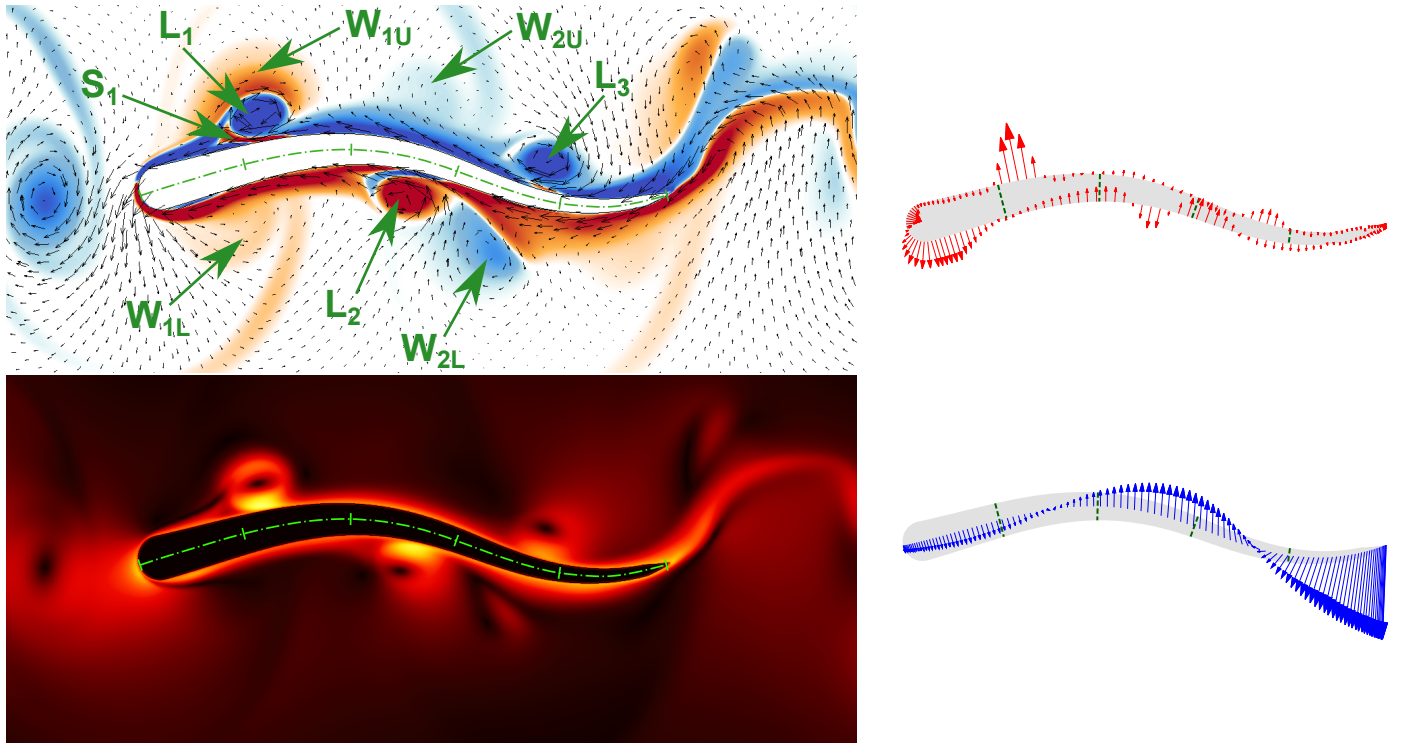
Figure 3: **Energetics data for a smart follower maximizing its swimming-efficiency.** (a) Swimming-efficiency, and (b) Cost of Transport for $IS_\eta$ (solid blue line) and $SS_\eta$ (dash-double-dot black line), normalized with respect to the CoT of a steady solitary swimmer. Four instances of maximum and minimum efficiency, which occur periodically throughout the simulation at times $(nT_p + 0.12)$, $(nT_p + 0.37)$, $(nT_p + 0.62)$, $(nT_p + 0.87)$, have been highlighted. $T_p = 1$ denotes the constant tail-beat period of the swimmers, whereas $n$ represents an integral multiple. The decline in $\eta$ at point $E$ ($t \approx 27.7$, $\eta = 0.86$) results from an erroneous manoeuvre at $t \approx 26.5$ (Supplementary Movie S4), which reveals the existence of a time-delay between actions and their consequences.

is used to regulate the undulations of two followers to maintain these target coordinates on either branch of the diverging wake, as shown in Fig. 1b and Supplementary Movie S1. The controlled motion yields an 11% increase in average swimming-efficiency for each of the followers (Fig. 5a), and a 5% reduction in each of their Cost of Transport. Overall, the group experiences a 7.4% increase in efficiency when compared to three isolated non-interacting swimers. The mechanism of energy-savings closely resembles that observed for the 2D swimmer; an oncoming wake-vortex ring (WR - Fig. 5b) interacts with the deforming body to generate a 'lifted-vortex' ring (LR - Fig. 5c). As this new ring proceeds along the length of the body, it modulates the follower's swimming-efficiency as observed in Fig. 5. Remarkably, the positioning of the lifted-ring at the instants of minimum and maximum swimming-efficiency resembles the corresponding positioning of lifted-vortices in the 2D case; a slight dip in efficiency corresponds to lifted-vortices interacting with the anterior section of the body (Fig. 5c and Supplementary Fig. S4), whereas an increase occurs upon their interaction with the midsection (Fig. 5d and Fig. 4).

These results showcase the remarkable capability of machine learning, and deep RL in particular, for discovering effective solutions that may not have been envisaged by humans, either owing to pre-existing biases, or due to the difficulty of anticipating the effects of delayed reactions by swimmers in complex flows. Finally, this study demonstrates that deep reinforcement learning can produce navigation algorithms for complex flow-fields, with promising implications for energy savings in autonomous robotic swarms.
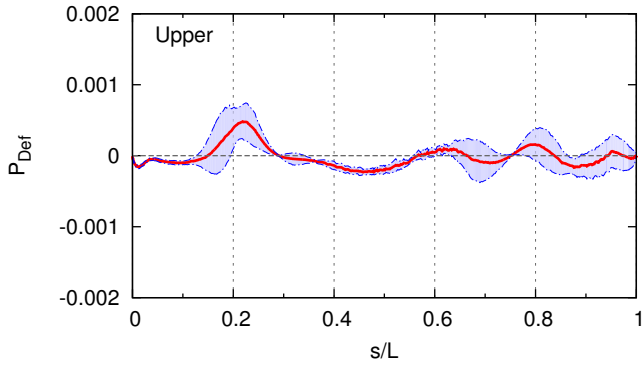
**Methods** We perform two-and three dimensional simulations of multiple self-propelled swimmers using wavelet adapted vortex methods[36] to discretise the velocity-vorticity form of the Navier-Stokes (NS) equations (in 2D), and their velocity-pressure form along with the pressure-projection[37] method (in 3D) using finite differences on a uniform computational grid. The body-geometry of the self-propelled swimmers is based on simplified models of a zebrafish. The swimmers adapt their motion using deep reinforcement learning. The learning process was greatly accelerated by employing recurrent neural networks with long-short term memory (RL-LSTM)[38] as a surrogate of the value function for the smart-swimmer. Additional details regarding the simulation methods and the reinforcement learning algorithm are provided in the Supporting Information.
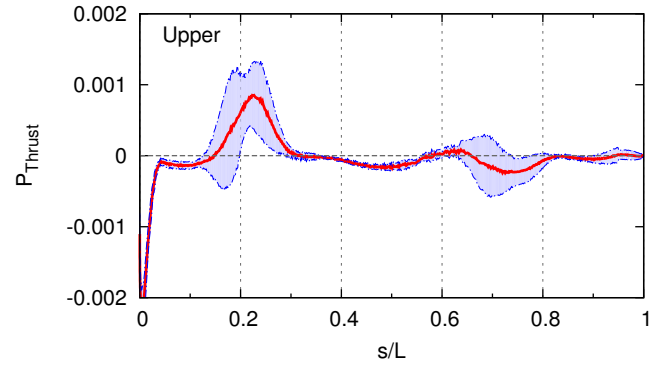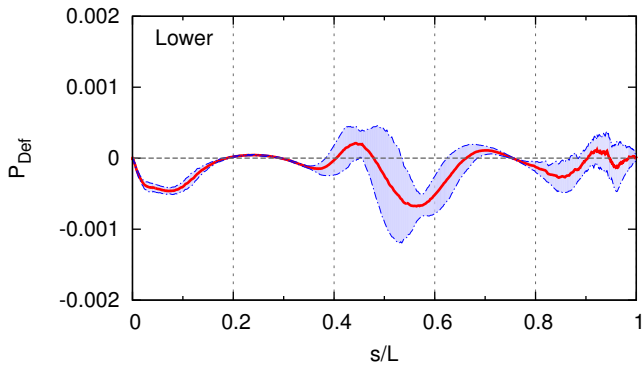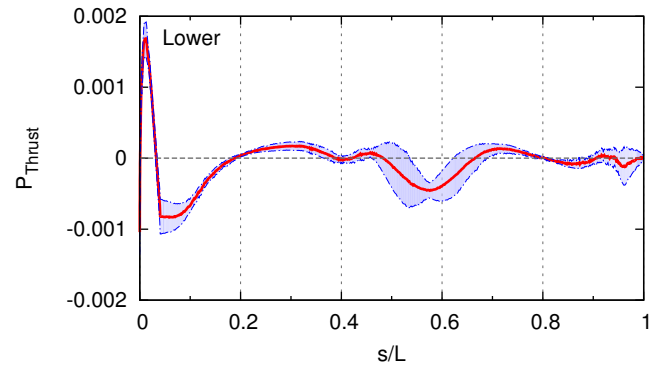
Figure 4: **Flow-field, and flow-induced forces for $IS_\eta$, corresponding to maximum efficiency.** (a) Vorticity field (red: positive, blue: negative) with velocity vectors shown as black arrows (top), and velocity magnitude shown in the lower panel (bright: high speed, dark: low speed). The snapshots correspond to $t = 26.12$, i.e., point $\eta_{max}(A)$ in Fig. 3a. In all the panels, demarcations are shown at every $0.2L$ along the body center-line for reference. The wake-vortices intercepted by the follower ($W_{1U}$, $W_{1L}$, $W_{2U}$, $W_{2L}$), the lifted-vortices created by interaction of the body with the flow ($L_1$, $L_2$, and $L_3$), and secondary-vorticity $S_1$ generated by $L_1$ have been annotated in the figure. (b) Flow-induced force-vectors (top) and body-deformation velocity (bottom) at $t = 26.12$. (c) Deformation-power, and (d) thrust-power (with negative values indicating drag-power) acting on the upper surface of follower. The red line indicates the average over 10 different snapshots ranging from $t = 30.12$ to $t = 39.12$. The envelope signifies the standard deviation among the 10 snapshots. (e) Deformation-power and (f) thrust-power on the lower (left-lateral) surface of the swimmer.
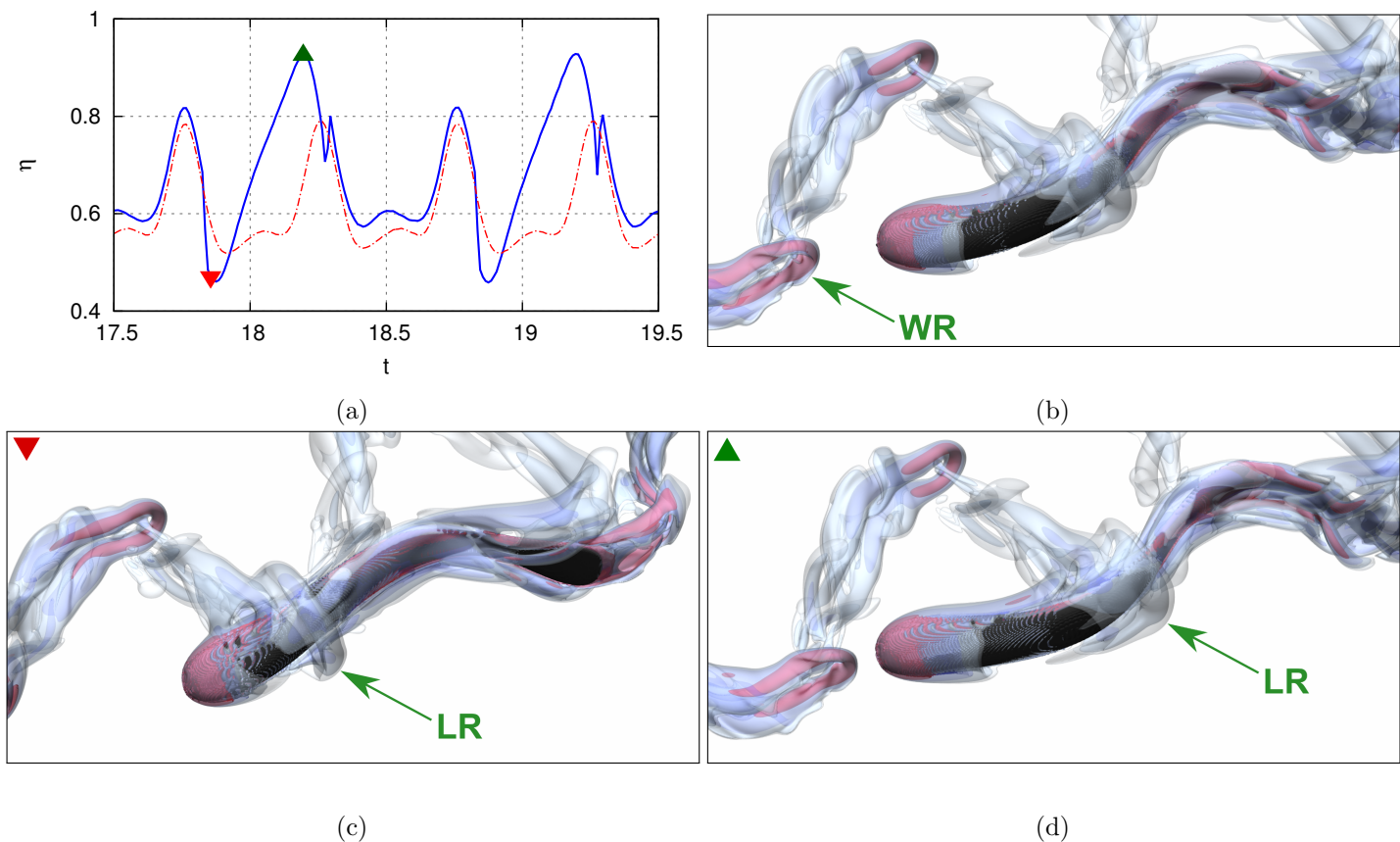
8

Figure 5: **3D swimmer interacting with wake vortex rings.** (a) Swimming-efficiency for a three-dimensional leader (dash-dot red line) and a follower (solid blue line) that adjusts its undulations via a Proportional-Integrator (PI) feedback controller to maintain a specified position in the wake. After an initial transient, the patterns visible in the efficiency-curves repeat periodically with $T_p$. Time-instances where the follower attains its minimum and maximum swimming-efficiency have been marked with an inverted red triangle, and an upright green triangle, respectively. The sudden jumps at $t \approx 18.3$ and $19.3$ correspond to adjustments made by the PI controller. (b) An oncoming wake-vortex ring (WR) is intercepted by the head of the follower, and generates a new 'lifted-vortex' ring (LR, panel c) similar to the 2D case (Fig. 4). As this ring interacts with the deforming body, it lowers the swimming-efficiency initially ($t \approx 17.8$: panels a and c), but provides a noticeable benefit further dowstream ($t \approx 18.2$, panels a and d).

9

# References

[1] Schmidt J (1923) Breeding places and migrations of the eel. *Nature* 111:51–54.

[2] Lang TG, Pryor K (1966) Hydrodynamic performance of porpoises (stenella attenuata). *Science* 152:531–533.

[3] Aleyev YG (1977) *Nekton.* (Springer Netherlands).

[4] Triantafyllou MS, Weymouth GD, Miao J (2016) Biomimetic Survival Hydrodynamics and Flow Sensing. *Annu. Rev. Fluid Mech.* 48:1–24.

[5] Breder CM (1965) Vortices and fish schools. *Zoologica-N.Y.* 50:97–114.

[6] Weihs D (1973) Hydromechanics of fish schooling. *Nature* 241:290–291.

[7] Shaw E (1978) Schooling Fishes: The school, a truly egalitarian form of organization in which all members of the group are alike in influence, offers substantial benefits to its participants. *Am. Sci.* 66:166–175.

[8] Pavlov DS, Kasumyan AO (2000) Patterns and mechanisms of schooling behavior in fish: A review. *J. Ichthyol.* 40:163–231.

[9] Burgerhout E, et al. (2013) Schooling reduces energy consumption in swimming male European eels, Anguilla anguilla L. *J. Exp. Mar. Biol. Ecol.* 448:66 – 71.

[10] Whittlesey RW, Liska S, Dabiri JO (2010) Fish schooling as a basis for vertical axis wind turbine farm design. *Bioinspir. Biomim.* 5(3):035005.

[11] Chapman JW, et al. (2011) Animal orientation strategies for movement in flows. *Curr. Biol.* 21:R861 – R870.

[12] Montgomery JC, Baker CF, Carton AG (1997) The lateral line can mediate rheotaxis in fish. *Nature* 389:960–963.

[13] Lyon EP (1904) On rheotropism. I. — Rheotropism in fishes. *Am. J. Physiol.* 12:149–161.

[14] Liao JC, Beal DN, Lauder GV, Triantafyllou MS (2003) Fish exploiting vortices decrease muscle activity. *Science* 302:1566–1569.

[15] Oteiza P, Odstrcil I, Lauder G, Portugues R, Engert F (2017) A novel mechanism for mechanosensory-based rheotaxis in larval zebrafish. *Nature* 547:445–448.

[16] Herskin J, Steffensen JF (1998) Energy savings in sea bass swimming in a school: measurements of tail beat frequency and oxygen consumption at different swimming speeds. *J. Fish Biol.* 53:366–376.

[17] Killen SS, Marras S, Steffensen JF, McKenzie DJ (2012) Aerobic capacity influences the spatial position of individuals within fish schools. *Proc. Biol. Sci.* 279:357–364.

[18] Ashraf I, et al. (2017) Simple phalanx pattern leads to energy saving in cohesive fish schooling. *Proc. Natl. Acad. Sci. U.S.A.*

[19] Pitcher TJ (1986) Functions of shoaling behaviour in teleosts in *The Behaviour of Teleost Fishes*, ed. Pitcher TJ. (Springer US, Boston, MA), pp. 294–337.

[20] Lopez U, Gautrais J, Couzin ID, Theraulaz G (2012) From behavioural analyses to models of collective motion in fish schools. *Interface Focus* 2:693–707.

[21] Daghooghi M, Borazjani I (2015) The hydrodynamic advantages of synchronized swimming in a rectangular pattern. *Bioinspir. Biomim.* 10:056018.

[22] Gazzola M, Hejazialhosseini B, Koumoutsakos P (2014) Reinforcement learning and wavelet adapted vortex methods for simulations of self-propelled swimmers. *SIAM J. Sci. Comput.* 36:B622–B639.

[23] Maertens AP, Gao A, Triantafyllou MS (2017) Optimal undulatory swimming for a single fish-like body and for a pair of interacting swimmers. *J. Fluid Mech* 813:301–345.

[24] Mnih V, , et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518:529–533.

[25] Müller UK, Smit J, Stamhuis EJ, Videler JJ (2001) How the body contributes to the wake in undulatory fish swimming. *J. Exp. Biol.* 204:2751–2762.

[26] Kern S, Koumoutsakos P (2006) Simulations of optimized anguilliform swimming. *J. Exp. Biol.* 209:4841–4857.

[27] Borazjani I, Sotiropoulos F (2008) Numerical investigation of the hydrodynamics of carangiform swimming in the transitional and inertial flow regimes. *J. Exp. Biol.* 211:1541–1558.

[28] Liao JC, Beal DN, Lauder GV, Triantafyllou MS (2003) The Kármán gait: novel body kinematics of rainbow trout swimming in a vortex street. *J. Exp. Biol.* 206:1059–1073.

[29] Sutton RS, Barto AG (1998) *Reinforcement learning: An introduction.* (MIT press, Cambridge, MA, USA).

[30] Gazzola M, Tchieu AA, Alexeev D, de Brauer A, Koumoutsakos P (2016) Learning to school in the presence of hydrodynamic interactions. *J. Fluid Mech.* 789:726–749.

[31] Reddy G, Celani A, Sejnowski TJ, Vergassola M (2016) Learning to soar in turbulent environments. *Proceedings of the National Academy of Sciences* 113(33):E4877–E4884.

[32] Colabrese S, Gustavsson K, Celani A, Biferale L (2017) Flow navigation by smart microswimmers via reinforcement learning. *Physical Review Letters* 118(15):158004–.

[33] Novati G, et al. (2017) Synchronisation through learning for two self-propelled swimmers. *Bioinspir. Biomim.* 12:036001.

[34] Weihs D (1975) *Swimming and Flying in Nature: Volume 2*, eds. Wu TYT, Brokaw CJ, Brennen C. (Springer US, Boston, MA), pp. 703–718.

[35] Bertsekas DP, Bertsekas DP, Bertsekas DP, Bertsekas DP (1995) *Dynamic programming and optimal control.* (Athena scientific Belmont, MA) Vol. 1.

[36] Rossinelli D, et al. (2015) MRAG-I2D: Multi-resolution adapted grids for remeshed vortex methods on multicore architectures. *J. Comput. Phys.* 288:1–18.

[37] Chorin AJ (1968) Numerical solution of the Navier-Stokes equations. *Math. Comp.* 22:745–762.

[38] Hochreiter S, Schmidhuber J (1997) Long short-term memory. *Neural Comput.* 9:1735–1780.

[39] Coquerelle M, Cottet GH (2008) A vortex level set method for the two-way coupling of an incompressible fluid with colliding rigid bodies. *J. Comput. Phys.* 227:9121–9137.

[40] Verma S, Abbati G, Novati G, Koumoutsakos P (2017) Computing the force distribution on the surface of complex, deforming geometries using vortex methods and brinkman penalization. *Int. J. Numer. Meth. Fluids* 85(8):484–501.

[41] Greengard L, Rokhlin V (1987) A fast algorithm for particle simulations. *J. Comput. Phys.* 73:325–348.

[42] Gholami A, Hill J, Malhotra D, Biros G (2015) AccFFT: A library for distributed-memory FFT on CPU and GPU architectures. *arXiv preprint arXiv:1506.07933.*

[43] Tytell ED, Lauder GV (2004) The hydrodynamics of eel swimming. *J. Exp. Biol.* 207:1825–1841.

[44] van Rees WM, Gazzola M, Koumoutsakos P (2013) Optimal shapes for anguilliform swimmers at intermediate reynolds numbers. *J. Fluid Mech.* 722:R3 1–12.

[45] Bellman RE (2010) *Dynamic Programming.* (Princeton University Press, Princeton, NJ, USA).

[46] van Hasselt H, Guez A, Silver D (2015) Deep reinforcement learning with double Q-learning. *CoRR, abs/1509.06461.*

[47] Mnih V, , et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518:529–533.

[48] Riedmiller M (2005) Neural fitted Q iteration – First experiences with a data efficient neural reinforcement learning method in *Machine Learning: ECML 2005: Lecture Notes in Computer Science, vol 3720*, eds. Gama J, Camacho R, Brazdil PB, Jorge AM, Torgo L. (Springer Berlin Heidelberg, Berlin, Heidelberg), pp. 317–328.

[49] Gers FA, Schmidhuber J, Cummins F (2000) Learning to forget: Continual prediction with LSTM. *Neural Comput.* 12(10):2451–2471.

[50] Lin LJ (1992) Ph.D. thesis (Carnegie Mellon University, Pittsburgh, PA, USA).

[51] Kingma DP, Ba J (2014) Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980 [cs.LG]*.

[52] Graves A, Schmidhuber J (2005) Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Netw.* 18:602–610.

[53] Hunt JCR, Wray AA, Moin P (1988) Eddies, streams, and convergence zones in turbulent flows in *Studying Turbulence Using Numerical Simulation Databases, 2. Report CTR-S88*. pp. 193–208.

# Supporting Information - Methods

**Simulation details.** The simulations presented here are based on the incompressible Navier-Stokes (NS) equations:

$$\nabla \cdot \boldsymbol{u} = 0 \tag{1}$$

$$\frac{\partial \boldsymbol{u}}{\partial t} + \boldsymbol{u} \cdot \nabla \boldsymbol{u} = -\frac{\nabla P}{\rho_f} + \nu \nabla^2 \boldsymbol{u} + \lambda \chi (\boldsymbol{u}_s - \boldsymbol{u}) \tag{2}$$

Each swimmer is represented on the computational grid via the characteristic function $\chi$, and interacts with the fluid by means of the penalty[39] term $\lambda \chi (\boldsymbol{u}_s - \boldsymbol{u})$, with $\lambda = 1e6$. $\boldsymbol{u}_s$ denotes the swimmer's combined translational, rotational, and deformation velocity, whereas $\boldsymbol{u}$ and $\nu$ correspond to the fluid velocity and viscosity, respectively. $P$ represents the pressure, and the fluid density is denoted by $\rho_f$.

The vorticity form of the NS equations was used for the two-dimensional simulations. A wavelet adaptive grid[36] with an effective resolution of $4096^2$ points was used to discretize a unit square domain. A lower effective resolution of $1024^2$ points was used for the training-simulations to minimize computational cost. The pressure-Poisson equation ($\nabla^2 P = -\rho_f \left( \nabla \boldsymbol{u}^T : \nabla \boldsymbol{u} \right) + \rho_f \lambda \nabla \cdot (\chi (\boldsymbol{u}_s - \boldsymbol{u}))$), necessary for estimating the distribution of flow-induced forces on the swimmers' bodies, was solved using the Fast Multipole Method.[40,41]

The three-dimensional simulations employed the pressure-projection method for solving the NS equations.[37] The simulations were parallelized via the CUBISM framework,[36] and used a uniform grid consisting of $2048 \times 1024 \times 256$ points in a domain of size $1 \times 0.5 \times 0.125$. The non-divergence-free deformation of the self-propelled swimmers was incorporated into the pressure-Poisson equation as follows:

$$\nabla^2 P = \frac{\rho_f}{\Delta t} \left( \nabla \cdot \boldsymbol{u}^\star - \chi \nabla \cdot \boldsymbol{u}_s \right), \tag{3}$$

where $\boldsymbol{u}^\star$ represents the intermediate velocity from the convection-diffusion-penalization fractional steps. Equation 3 was solved using a distributed Fast Fourier Transform library (AccFFT[42]).

**Flow-induced forces, and energetics variables.** The pressure-induced and viscous forces acting on the swimmers are computed as follows:[40]

$$\boldsymbol{dF}_P = -P\boldsymbol{n} \, dS \tag{4}$$

$$\boldsymbol{dF}_\nu = 2\mu \boldsymbol{D} \cdot \boldsymbol{n} \, dS \tag{5}$$

Here, $P$ represents the pressure acting on the swimmer's surface, $\boldsymbol{D} = \left( \nabla \boldsymbol{u} + \nabla \boldsymbol{u}^T \right)/2$ is the strain-rate tensor on the surface, and $dS$ denotes the infinitesimal surface area. Since self-propelled swimmers generate zero net average thrust (and drag) during steady swimming, we determine the instantaneous thrust as follows:

$$\text{Thrust} = \frac{1}{2\|\boldsymbol{u}\|} \iint \left( \boldsymbol{u} \cdot \boldsymbol{dF} + |\boldsymbol{u} \cdot \boldsymbol{dF}| \right), \tag{6}$$

where $\boldsymbol{dF} = \boldsymbol{dF}_P + \boldsymbol{dF}_\nu$. Similarly, the instantaneous drag may be determined as:

$$\text{Drag} = \frac{1}{2\|\boldsymbol{u}\|} \iint \left( \boldsymbol{u} \cdot \boldsymbol{dF} - |\boldsymbol{u} \cdot \boldsymbol{dF}| \right) \tag{7}$$

Using these quantities, the thrust-, drag-, and deformation-power are computed as:

$$P_{Thrust} = \text{Thrust} \cdot \|\boldsymbol{u}\| \tag{8}$$

$$P_{Drag} = -\text{Drag} \cdot \|\boldsymbol{u}\| \tag{9}$$

$$P_{Def} = -\iint \boldsymbol{u}_{Def} \cdot \boldsymbol{dF} \tag{10}$$

where $\boldsymbol{u}_{Def}$ represents the deformation-velocity of the swimmer's body. The double-integrals in these equations represent surface-integration over the swimmer's body, and yield measurements for time-series analysis. On the

other hand, only the integrand is evaluated when surface-distributions of thrust-, drag-, or deformation-power are required (as in Figs. 4c to 4f).

The instantaneous swimming-efficiency is based on a modified form of the Froude efficiency proposed in ref.:[43]

$$\eta = \frac{P_{Thrust}}{P_{Thrust} + \max(P_{Def}, 0)} \tag{11}$$

To compute both $\eta$ and the Cost of Transport (CoT), we neglect negative values of $P_{Def}$, which can result from beneficial interactions of the smart-swimmer with the leader's wake:

$$CoT(t) = \frac{\int_{t-T_p}^{t} \max(P_{Def}, 0) dt}{\int_{t-T_p}^{t} \|\boldsymbol{u}\| dt} \tag{12}$$

This restriction accounts for the fact that the elastically rigid swimmer may not store energy furnished by the flow, and yields a conservative estimate of potential savings in the CoT. We note that percentage-changes in $P_{Def}$, reported in the main text and the supplementary section, have been computed using this bounded value to avoid overstating any potential benefits.

**Swimmer shape and kinematics.** The Reynolds number of the self-propelled swimmers is computed as $Re = L^2/(\nu T_p)$. The body-geometry is based on a simplified model of a zebrafish.[44] The half-width of the 2D profile is described as follows:

$$w(s) = \begin{cases} \sqrt{2w_h s - s^2} & 0 \leq s < s_b \\ w_h - (w_h - w_t)\left(\dfrac{s - s_b}{s_t - s_b}\right) & s_b \leq s < s_t \\ w_t \dfrac{L - s}{L - s_t} & s_t \leq s \leq L \end{cases} \tag{13}$$

where $s$ is the arc-length along the midline of the geometry, $L = 0.1$ is the body length, $w_h = s_b = 0.04L$, $s_t = 0.95L$, and $w_t = 0.01L$. For 3D simulations, the geometry is comprised of elliptical cross sections, with the half-width $w(s)$ and half-height $h(s)$ described via cubic B-splines.[44] Six control-points define the half-width: $(s/L, w/L) = [(0.0, 0.0), (0.0, 0.089), (1/3, 0.017), (2/3, 0.016), (1.0, 0.013), (1.0, 0.0)]$; whereas eight control-points define the half-height: $(s/L, h/L) = [(0.0, 0.0), (0.0, 0.055), (0.2, 0.068), (0.4, 0.076), (0.6, 0.064), (0.8, 0.0072), (1.0, 0.11), (1.0, 0.0)]$. The length was set to $L = 0.2$, which keeps the grid-resolution, i.e., the number of points along the fish midline, comparable to the 2D simulations. Body-undulations for both 2D and 3D simulations were generated as a travelling-wave defining the curvature along the midline:

$$k(s, t) = A(s) \sin\left(\frac{2\pi t}{T_p} - \frac{2\pi s}{L}\right) \tag{14}$$

Here $A(s)$ is the curvature amplitude and varies linearly from $A(0) = 0.82$ to $A(L) = 5.7$.

**Reinforcement Learning.** Reinforcement learning (RL)[29] is a process by which an agent (in this case, the smart-swimmer) learns to earn rewards through trial-and-error interaction with its environment. At each turn, the agent observes the state of the environment $s_n$ and performs an action $a_n$, which influences both the transition to the next state $s_{n+1}$ and the reward received $r_{n+1}$. The agent's goal is to learn the optimal control policy $a_n = \pi^*(s_n)$ which maximises the action value $Q^*(s_n, a_n)$, defined as the sum of discounted future rewards:

$$Q^*(s_n, a_n) = \max_\pi \mathbb{E}\left(r_{n+1} + \gamma r_{n+2} + \gamma^2 r_{n+3} + \ldots \mid a_m = \pi(s_m) \; \forall m \in [n+1, \mathcal{T}]\right) \tag{15}$$

Here, $\mathcal{T}$ denotes the terminal state of a training-simulation, and the discount factor $\gamma$ is set to 0.9. The optimal action-value function $Q^*(s_n, a_n)$ is a fixed point of the Bellman equation: $Q^*(s_n, a_n) = \mathbb{E}\left[r_{n+1} + \gamma \max_{a'} Q^*(s_{n+1}, a')\right]$.[45] We approximate $Q^*(s_n, a_n)$ using a neural network[46–48] with weights $w_k$, which are updated iteratively to minimize the temporal difference error:

$$\text{TD}_{\text{err}} = \mathbb{E}_{s_n, a_n, s_{n+1}}\left[r_{n+1} + \gamma Q(s_{n+1}, a'; \text{w}_-) - Q(s_n, a_n; \text{w}_k)\right] \tag{16}$$

14

Here, $w_-$ is a set of target weights, and $a'$ is the best action in state $s_{n+1}$ computed with the current weights ($a' = \arg\max_a Q(s_{n+1}, a; w_k)$). The target weights $w_-$ are updated towards the current weights as $w_- \leftarrow (1 - \alpha)w_- + \alpha w_k$, where $\alpha = 10^{-4}$ is an under-relaxation factor used to stabilize the algorithm.[47]

**States and actions.** The six observed-state variables perceived by the learning agent include $\Delta x$, $\Delta y$, $\theta$, the two most recent actions taken by the agent, and the current tail-beat 'stage' $\mod(t, T_p)/T_p$. The permissible range of the observed-state variables is limited to: $1 \leq \Delta x/L \leq 3$; $|\Delta y|/L \leq 1$ (boundary depicted by $R_{end}$ in Supplementary Fig. S7); and $|\theta| \leq \pi/2$. If the agent exceeds any of these thresholds, the training-simulation terminates and the agent receives a terminal reward $R_{end} = -1$.

The smart-swimmer (or agent) is capable of manoeuvering by actively manipulating the curvature-wave travelling down the body. This is accomplished by linearly superimposing a piecewise function on the baseline curvature $k(s,t)$ (equation 14):

$$k_{\text{Agent}}(s,t) = k(s,t) + A(s)M(t, T_p, s, L) \tag{17}$$

The curve $M(t, T_p, s, L)$ is composed of 3 distinct segments:

$$M(t, T_p, s, L) = \sum_{j=0}^{2} b_{n-j} \cdot m\left(\frac{t - t_{n-j}}{T_p} - \frac{s}{L}\right) \tag{18}$$

The curve $m$ is a clamped cubic spline with $m(0) = m'(0) = 0$, $m(1/2) = m'(1/2) = 0$, and $m(1/4) = 1$, $m'(1/4) = 0$. $t_n$ represents the time-instance when action $a_n$ is taken, whereas $b_n$ represents the corresponding control-amplitude, which may take five discrete values: 0, $\pm 0.25$, and $\pm 0.5$.

**Neural network architecture.** One of the assumptions in RL is that the transition probability to a new state $s_{n+1}$ is independent of the previous transitions, given $s_n$ and $a_n$, i.e.,:

$$p(s_{n+1} \,|\, s_n, a_n) = p(s_{n+1} \,|\, s_n, a_n, \ldots, s_0, a_0) \tag{19}$$

This assumption is invalidated whenever the agent has a limited perception of the environment. In most realistic cases the agent receives an observation $o_n$ rather than the complete state of the environment $s_n$. Therefore, past observations carry information relevant for future transitions (i.e., $p(o_{n+1} \,|\, o_n, a_n) \neq p(o_{n+1} \,|\, o_n, a_n, \ldots, o_0, a_0)$), and should be taken into account in order to make optimal decisions. This operation can be approximated by a Recurrent Neural Network (RNN), which can learn to compute and remember important features in past observations. In this work we approximate the action-value function with a LSTM-RNN[49] composed of three layers of 24 fully connected LSTM cells each, and terminating in a linear layer (Supplementary Fig. S3). The last layer computes a vector of action-values $\mathbf{q}_n = Q(o_n; y_{n-1}, w_k)$ with one component $q_n^{(a)}$ for each possible action $a$ available to the agent ($y_{n-1}$ represents the activation of the network at the previous turn).

**Training procedure.** During training, both the leader and the follower (learning agent) start from rest. The leader swims steadily along a straight line, whereas the follower manoeuvers according to the actions supplied to it. Multiple independent simulations run simultaneously, with each of these sending the current observed-state $o_n$ of the agent to a central processor, and in turn receiving the next action $a_n$ to be performed. The central processor computes $a_n$ using an $\epsilon$-greedy policy (with $\epsilon$ gradually annealed from 1 to 0.1) from the most recently updated $Q$ function. Once a training-simulation reaches a terminal state (e.g., the follower hits the boundary labelled $R_{end}$ in Supplementary Fig. S7), all the messages exchanged between the simulation and the central processor are appended to a training set of sequences $\mathcal{R}$.[50] In the meantime, the network is continually updated by sampling $B$ sequences from the set $\mathcal{R}$, according to algorithm 1. The batch gradient $\Delta w$ is computed with back propagation through time (BPTT).[52] The network weights are then updated with the Adam stochastic optimization algorithm.[51]

**Algorithm 1: Asynchronous recurrent DQN algorithm.**

initialize network $w_0$ and target network $w_- = w_0$;
initialize set of transition sequences $\mathcal{R} = \emptyset$;
**repeat**
    $N \leftarrow 0$;
    sample batch of $B$ sequences from $\mathcal{R}$;
    **for** *sequence* $j \in [1, \ldots, B]$ **do**
        $[\mathbf{q}_{j,0}, y_{j,0}] = Q(o_{j,0}; \emptyset, w_k)$;
        **for** *turns* $n \in [0, \ldots, \mathcal{T}_j]$ **do**
            $[\mathbf{q}_{j,n+1}, y_{j,n+1}] = Q(o_{j,n+1}; y_{j,n}, w_k)$;
            $[\tilde{\mathbf{q}}_{j,n+1}, \tilde{y}_{j,n+1}] = Q(o_{j,n+1}; y_{j,n}, w_-)$;
            $a' = \arg\max_a \left[ q_{j,n+1}^{(a)} \right]$;
            **if** $s_{j,n+1}$ *is terminal* **then**
                $e_{j,n} = r_{j,n+1} - q_{j,n}^{(a_n)}$;
            **else**
                $e_{j,n} = r_{j,n+1} + \gamma \tilde{q}_{j,n+1}^{(a')} - q_{j,n}^{(a_n)}$;
            **end**
            $N \leftarrow N + 1$;
        **end**
    **end**
    perform BPTT: $\Delta w = \frac{1}{N} \sum_j \sum_n e_{j,n} \nabla_w q_{j,n}^{(a_n)}$;
    update weights $w_{k+1}$ with Adam algorithm[51];
    update target network: $w_- \leftarrow (1 - \alpha)w_- + \alpha w_{k+1}$;
    $k \leftarrow k + 1$;
**until** $Q(o, a; w_k) = Q^*(o, a)$;

**Proportional-Integral feedback controller.** The PI controller modulates the 3D follower's body-kinematics, which allows it to maintain a specific position $(x_{tgt}, y_{tgt}, z_{tgt})$ relative to the leader:

$$k(s, t) = \alpha(t)A(s) \left[ \sin\left( \frac{2\pi t}{T_p} - \frac{2\pi s}{L} \right) + \beta(t) \right] \tag{20}$$

The factor $\alpha(t)$ modifies the undulation envelope, and controls the acceleration or deceleration of the follower based on its streamwise distance from the target position:

$$\alpha(t) = 1 + f_1 \left( \frac{x - x_{tgt}}{L} \right) \tag{21}$$

The term $\beta(t)$ adds a baseline curvature to the follower's midline to correct for lateral deviations:

$$\beta(t) = \frac{y_{tgt} - y}{L} \left( f_2|\theta| + f_3|\hat{\theta}| \right) \tag{22}$$

Here, $\theta$ represents the follower's yaw angle about the $z$-axis, and $\hat{\theta}$ is its exponential moving average: $\hat{\theta}_{t+1} = \frac{1 - \Delta t}{T_p} \hat{\theta}_t + \frac{\Delta t}{T_p} \theta$. The swimmers' $z$-positions remain fixed at $z_{tgt}$, as out-of-plane motion is not permitted. The controller-coefficients were selected to have a minimal impact on regular swimming kinematics, which allows for a direct comparison of the follower's efficiency to that of the leader:

$$f_1 = 1 \tag{23}$$

$$f_2 = \max(0, 50 \operatorname{sign}(\theta \cdot (y_{tgt} - y))) \tag{24}$$

$$f_3 = \max(0, 20 \operatorname{sign}(\hat{\theta} \cdot (y_{tgt} - y))) \tag{25}$$

# Supporting Information - Supplementary Text, Figures, and Movies

**Body-deformation during autonomous manoeuvres.** The extent of body-bending that swimmers $IS_\eta$ and $IS_d$ undergo when manoeuvring is compared quantitatively in Supplementary Fig. S1. A qualitative comparison was presented in Fig. 2d. We observe that the body-deformation of $IS_d$ is noticeably higher than that of a steady swimmer (with relative curvature 1), which implies a tendency to take aggressive turns. The deformation for swimmer $IS_\eta$ is markedly lower, which plays an instrumental role in reducing the power required for undulating the body against flow-induced forces.

**Comparison of four different swimmers.** The performance metrics for four different swimmers are compared in Supplementary Fig. S2. Interacting swimmer $IS_d$ occasionally attains higher speed than $IS_\eta$ (Supplementary Fig. S2a), but at the cost of much higher energy expenditure (Supplementary Fig. S2c and Table 1). Moreover, the

| | $IS_\eta$ | $SS_\eta$ | $IS_d$ | $SS_d$ |
|---|---|---|---|---|
| $\eta$ | 1.0 | 0.76 | 0.77 | 0.66 |
| CoT | 1.0 | 1.56 | 3.96 | 3.86 |
| $P_{Def}$ | 1.0 | 1.41 | 3.90 | 3.28 |
| $P_{Thrust}$ | 1.0 | 0.66 | 2.33 | 1.48 |

Table 1: **Comparison of energetics metrics for the four swimmers.** Averaged values computed for the data shown in Supplementary Fig. S2. All the values shown have been normalized with respect to the corresponding value for $IS_\eta$.

speeds of solitary swimmers $SS_\eta$ and $SS_d$ are lower than those of either interacting swimmer ($IS_\eta$ and $IS_d$), which suggests that wake-interactions may benefit a follower regardless of the goal being pursued. In Supplementary Fig. S2d $P_{Def}$ attains negative values only for $IS_\eta$, which is indicative of maximum benefit extracted from flow-induced forces. Both $IS_d$ and $SS_d$ are capable of generating significantly higher thrust-power than $IS_\eta$, but suffer from larger deformation-power, and consequently, lower swimming-efficiency. Comparing the columns for $IS_\eta$ and $SS_\eta$ in Table S1, we note that interacting with a preceding wake has a measurable impact on swimming-performance; $IS_\eta$ is approximately 32% more efficient than $SS_\eta$, spends 36% less energy per unit distance travelled, requires 29% less power for body-undulations, and generates 52% higher thrust-power. Wake-interactions yield energetics benefits even for the swimmer actively minimizing lateral displacement from the leader, primarily by increasing thrust-power, as can be surmised by comparing the data for $IS_d$ and $SS_d$ in Supplementary Table 1.

**Uncovering underlying time-dependencies.** While it is relatively straightforward to maintain a particular tandem formation via feedback control (when the follower strays too far to one side, a feedback controller can relay instructions to veer in the opposite direction), the same is not true for maximizing swimming-efficiency. It is difficult to formulate a simple set of a-priori rules for maximizing efficiency, especially in dynamically evolving conditions. This happens because: 1) the swimmer perceives only a limited representation of its environment (Fig. 1d); and 2) there may be measurable delay between an action and its impact on the reward received over the long term. These traits make deep RL ideal for determining the optimal policy when maximizing swimming-efficiency, especially when augmented with recurrent neural networks (Supplementary Fig. S3). These network architectures are adept at discovering and exploiting long-term time-dependencies.

**Flow-interactions at the instant of minimum swimming-efficiency.** The instant when swimmer $IS_\eta$ attains the lowest efficiency during each half-period ($\eta_{min}(D)$ in Fig. 3a) is examined in Supplementary Fig. S4. The mean $P_{Def}$ curve is mostly positive on both the lower and upper surfaces, with large positive peaks generated by interaction with the wake- and lifted-vortices. This increase in effort is not offset sufficiently by an increase in $P_{Thrust}$, resulting in low swimming-efficiency. Compared to the instance of maximum efficiency (Fig. 4), increased effort is required in the head region, along with an increase in thrust-production by the tail section $s > 0.7L$.
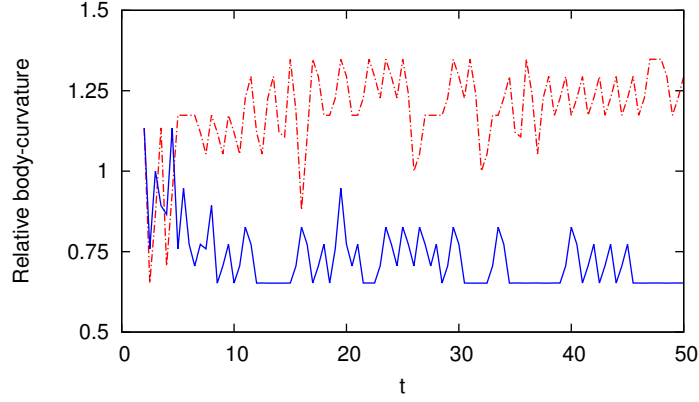
Figure S 1: **Midline curvature.** Severity of body-deformation for the swimmers $IS_\eta$ (solid blue line) and $IS_d$ (dash-dot red line), shown for 50 tail-beat periods starting from rest. The relative body-curvature is computed as $\Sigma_{i=1}^{6}|\kappa_i|$, normalized with the same metric for a solitary swimmer executing steady motion ($\kappa_i$ represents the curvature at 6 control points along a swimmer's body).

**Slight deviations impact performance.**    To examine the impact of small deviations in $IS_\eta$'s trajectory on its performance, we compare two different time-instances (at the same tail-beat stage) in Supplementary Fig. S5. At $t \approx 26.5$, $IS_\eta$ deviates slightly to the left of its steady trajectory (Supplementary Movie S4), which throws it out of synchronization with the oncoming wake-vortices. The resulting reduction in efficiency at $t \approx 27.5$ indicates that even slight deviations are capable of impacting performance, and that there may be a measurable delay between actions and consequences. However, the smart-swimmer autonomously corrects for such deviations, and is able to quickly recover its optimal behaviour.

**Correlation with the flow-field**    The correlation-coefficient curve shown in Fig. 2c, and the correlation map shown in Supplementary Fig. S6, were computed as follows:

$$\rho(\boldsymbol{u}, \boldsymbol{u}_{\text{head}}) = \frac{\text{cov}\left(\boldsymbol{u}(x,y), \boldsymbol{u}_{head}\right)}{\sigma_{\boldsymbol{u}(x,y)}\ \sigma_{\boldsymbol{u}_{\text{head}}}} = \frac{\sum_t \boldsymbol{u}(x,y,t) \cdot \boldsymbol{u}_{\text{head}}(t)}{\sqrt{\sum_t \|\boldsymbol{u}(x,y,t)\|^2}\sqrt{\sum_t \|\boldsymbol{u}_{\text{head}}(t)\|^2}} \tag{26}$$

Here, $\boldsymbol{u}(x,y,t)$ was recorded in the wake of a solitary swimmer, whereas $\boldsymbol{u}_{\text{head}}(t)$ was recorded at the swimmer's head. Maxima in $\rho(\boldsymbol{u}, \boldsymbol{u}_{\text{head}})$ provide an estimate for the coordinates where a follower's head-movements would exhibit long-term synchronization with an undisturbed wake.

**Limiting the exploration space.**    During training, the range of values that a smart-follower's states can take are constrained, as mentioned previously. This prevents excessive exploration of regions that involve no wake-interactions, and helps to minimize the computational cost of training-simulations. The limits of the bounding box (shown in Supplementary Fig. S7) are kept sufficiently large to provide the follower ample room to swim clear of the unsteady wake, if it determines that interacting with the wake is unfavourable.

**Power distribution in the presence/absence of a preceding wake.**    To determine the extent to which wake-induced interactions alter the distribution of $P_{Def}$ and $P_{Thrust}$, both of which influence overall swimming-efficiency, we compare these quantities for $IS_\eta$ and $SS_\eta$ in Supplementary Fig. S8. A similar comparison for $IS_d$ and $SS_d$ is shown in Supplementary Fig. S9. For $IS_\eta$, a greater variation in $P_{Def}$ and $P_{Thrust}$ is observed (broad envelopes in Supplementary Figs. S8a and S8b), compared to the solitary swimmer $SS_\eta$ (Supplementary Figs. S8c and S8d). This is caused by $IS_\eta$'s interactions with the unsteady wake, which is absent for $SS_\eta$. The average $P_{Def}$ for $IS_\eta$ shows distinct negative troughs near the head ($s/L < 0.2$, Supplementary Fig. S8a) and at $s/L = 0.6$. A lack of similar troughs for $SS_\eta$ (Supplementary Fig. S8c) implies that these benefits originate exclusively from wake-induced interactions. There is no apparent difference in drag for both $IS_\eta$ and $SS_\eta$ in the pressure-dominated
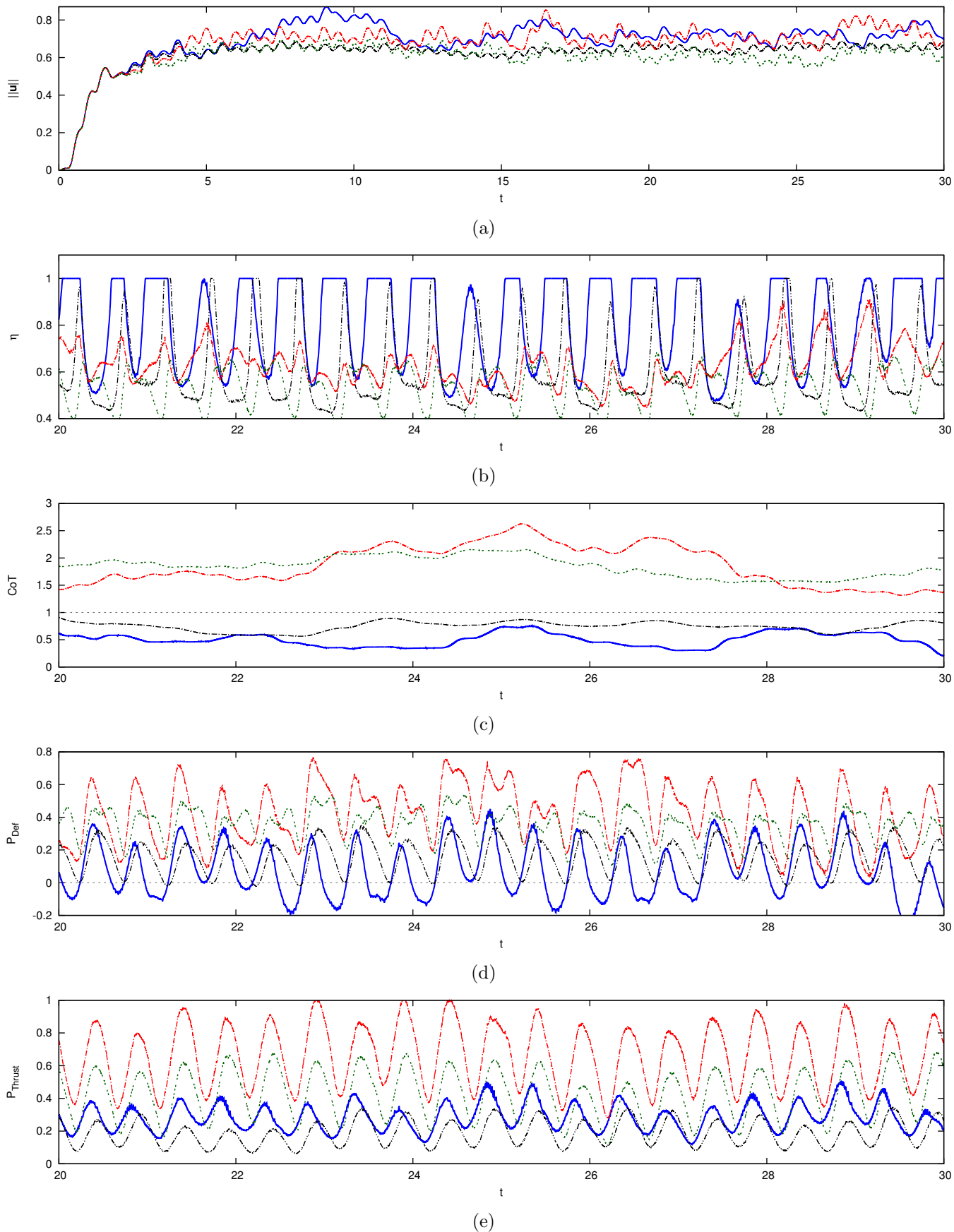
Figure S 2: **Performance metrics for four different swimmers.** Plots comparing (a) speed, (b) $\eta$, (c) CoT, (d) deformation-power , and (e) thrust-power for four different swimmers. The solid blue line corresponds to swimmer $IS_\eta$, the dash-double-dot black line to swimmer $SS_\eta$ (a solitary swimmer executing actions identical to $IS_\eta$), the dash-dot red line to swimmer $IS_d$, and the double-dot green line to swimmer $SS_d$ (a solitary swimmer executing actions identical to $IS_d$).

19

Figure S 3: **Schematic of the Recurrent Neural Network (RNN).** The RNN used in this study is composed of 3 LSTM layers, consisting of 24 cells (green blocks) each. The input layer (pink block) of the network comprises the 6 observed-state variables. The black arrows between different layers indicate all-to-all connections. The purple arrows indicate recurrent connections within each LSTM layer. The last layer consists of 5 output neurons (orange) with linear activation.
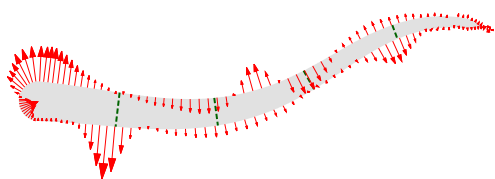
Figure S 4: **Flow-field and flow-induced forces for $IS_\eta$, corresponding to minimum efficiency.** (a) Vorticity field with the velocity vectors shown (top), and velocity magnitude (bottom) at $t = 26.87$ (point $\eta_{min}(D)$ in Fig. 3). (b) Flow-induced force-vectors (top) and body-deformation velocity (bottom) at this instance. (c,d) Deformation-power and thrust-power acting on the upper (right lateral) surface of follower. The red line indicates the average over 10 different snapshots ranging from $t = 30.87$ to $t = 39.87$. The envelope denotes the standard deviation among the 10 snapshots. (e,f) Deformation-power and thrust-power on the lower (left lateral) surface of the fish.
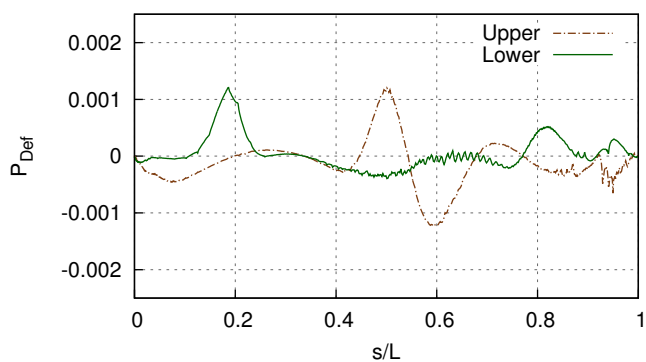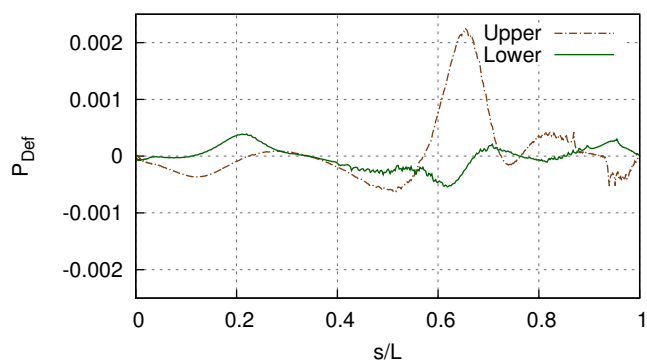
(a)  (b)

(c)  (d)

(e)  (f)

Figure S 5: **Deviations impact performance.** Comparison of two instances when a maximum in the swimming-efficiency is expected. The deformed shape and deformation-velocity for the two instances are similar, but differences in the flow-field influence efficiency. Panels on the left hand side of the page show data for $IS_\eta$ at $t \approx 33.7$ ($\eta = 1$), whereas those on the right hand side correspond to $t \approx 27.7$ ($\eta = 0.86$). (a, b) Vorticity, velocity vectors, and velocity magnitude at the two time instances. A slight deviation in the follower's approach to the wake causes a noticeable change in the surrounding vortices, as well as in the velocity induced near the surface. The regions highlighting differences have been marked as $R_1$, $R_2$, $R_3$, and $R_4$. (c, d) A comparison of the surface force-vectors and body-deformation velocity. (e,f) There are notable differences in the distribution of $P_{Def}$ on the upper and lower surfaces.
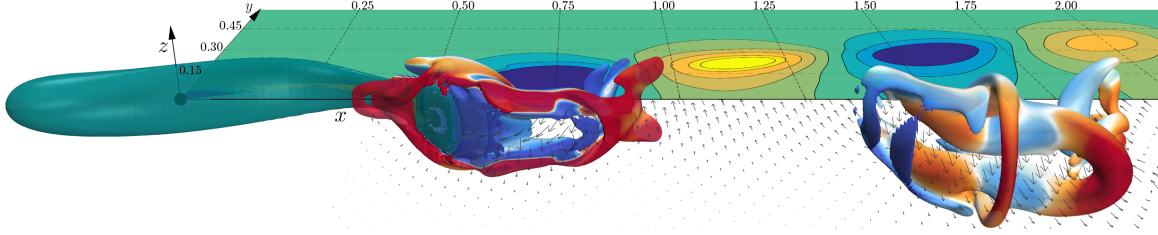
Figure S 6: **Correlation map.** The horizontal plane on the right side of the swimmer depicts the correlation-coefficient described by Equation 26. Areas of high correlation are denoted as yellow regions, whereas those of low correlation are shown in blue. The vortex rings shed are shown on the swimmer's left side, along with the velocity vectors on the left horizontal plane.
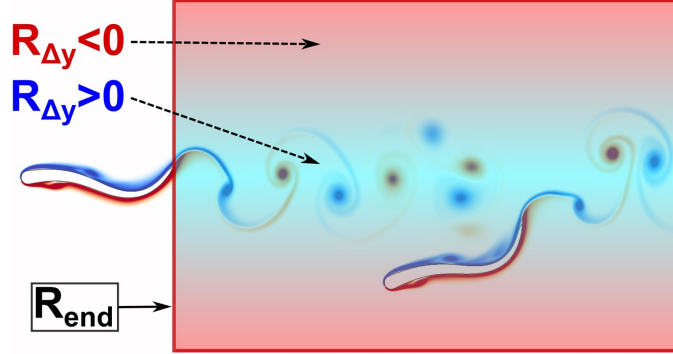


Figure S 7: **Reward for** $IS_d$**.** Visual representation of reward assigned to smart-swimmer $IS_d$, whose goal is to minimize its lateral displacement from the leader.

region close to the head ($s \approx 0$). However, wake-induced interactions provide a pronounced increase in thrust-power generated by the midsection for $IS_\eta$ (compare Supplementary Figs. S8b and S8d, $0.2 < s/L < 0.4$). Among all of the four swimmers compared, only $IS_\eta$ shows a distinct negative $P_{Def}$ region close to the head ($s < 0.2L$), which further supports the occurrence of head-motion synchronization with flow-induced forces, when efficiency is maximized. Comparing the deformation- and thrust-power distribution for $IS_d$ and $SS_d$ in Supplementary Fig. S9 provides additional evidence that wake-interactions have a marked impact on swimming-energetics.

**Supplementary Movie S1.** 3D simulation of three nonautonomous swimmers, in which the leader swims steadily, and the two followers maintain specified relative positions such that they interact favourably with the leader's wake. The flow-structures have been visualized using isosurfaces of the Q-criterion.[53]

**Supplementary Movie S2.** 2D simulation of a pair of swimmers, in which the leader swims steadily, and the follower ($IS_\eta$) takes autonomous decisions to interact favourably with the wake. The upper panel (labelled 'ω') shows the vorticity field generated by the swimmers, whereas the second panel (labelled 'v') shows the lateral flow-velocity. The smart-swimmer appears to synchronize the motion of its head with the lateral flow-velocity, which allows it to increase its swimming-efficiency. The lower panels show the energetics metrics, namely, the swimming efficiency $\eta$, the thrust-power $P_{Thrust}$, the deformation-power $P_{Def}$, and the Cost of Transport (CoT).

**Supplementary Movie S3.** 2D simulation of a pair of swimmers, where the leader performs random actions, and the follower takes autonomous decisions to benefit from the flow-field. The smart-follower, which was trained with a steadily-swimming leader, is able to adapt to the erratic leader's behaviour without any further training. Remarkably, the follower chooses to interact deliberately with the wake in order to maximize its long-term swimming-efficiency, even though it has the option to swim clear of the unsteady flow-field.
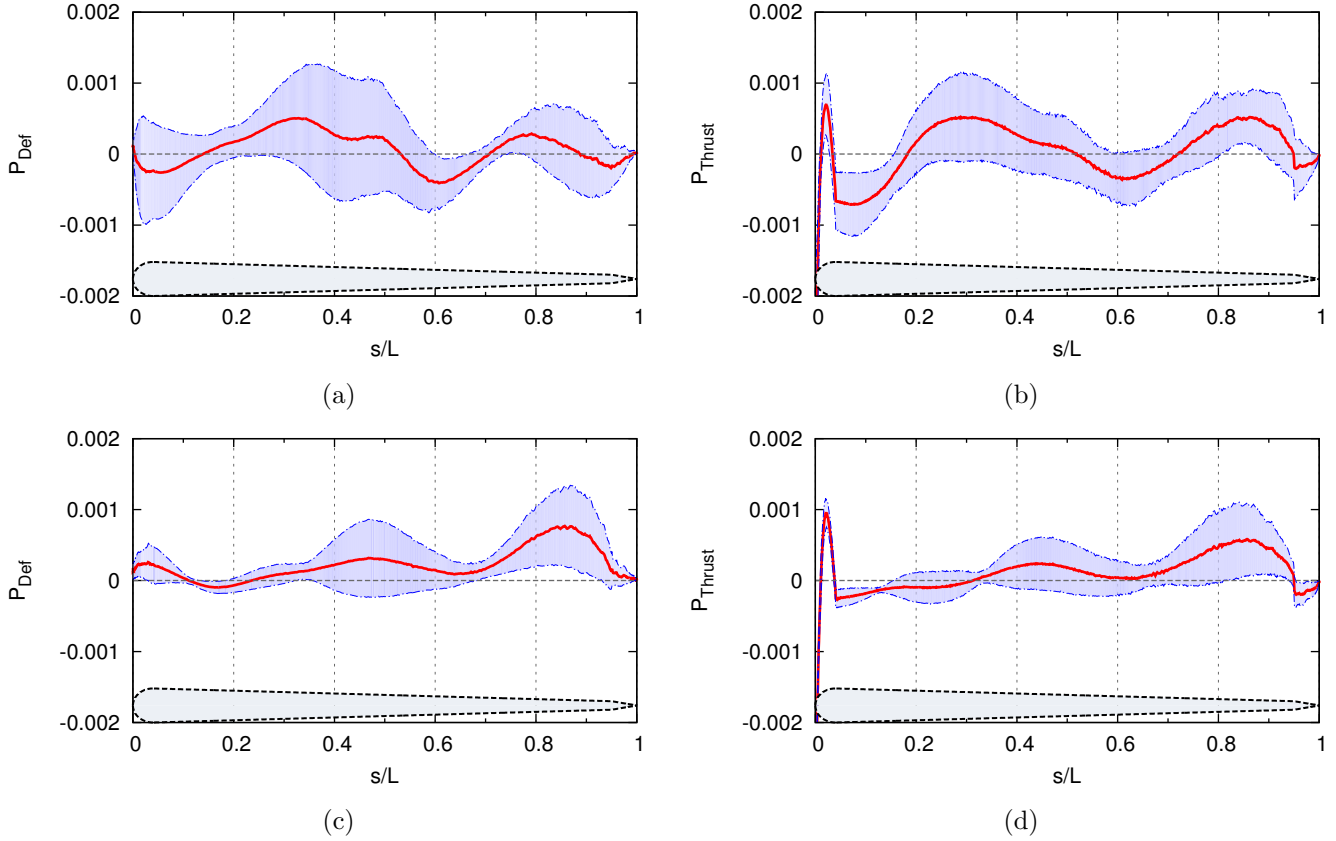
Figure S 8: **Power distribution.** Deformation-power and thrust-power distribution along the body of (a,b) swimmer $IS_\eta$, and (c,d) swimmer $SS_\eta$. The solid red line indicates the average over a single tail-beat period (from $t = 26$ to $t = 27$), whereas the envelope denotes the standard-deviation. The silhouettes at the bottom of each panel represent the fish body.
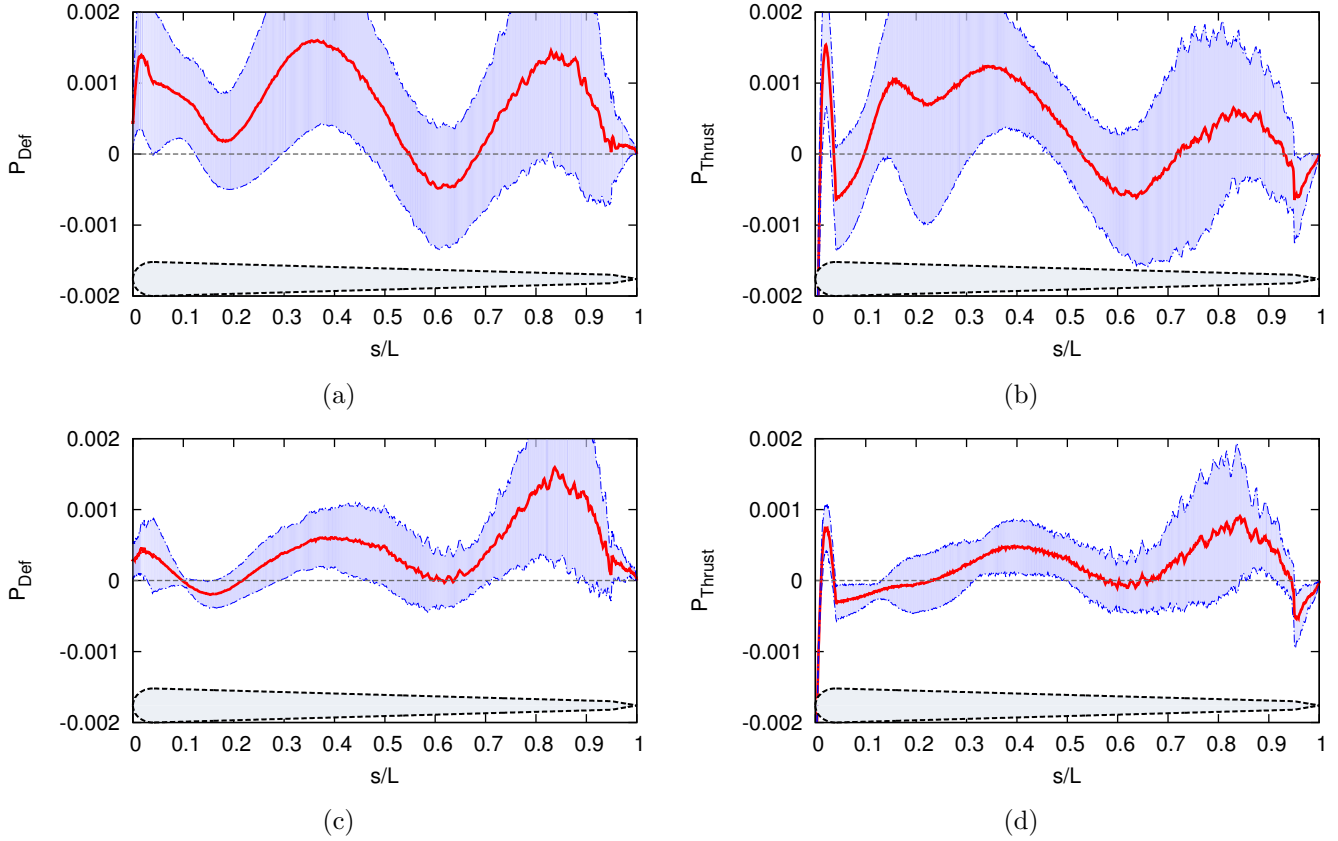
(a)

(b)

(c)

(d)

Figure S 9: **Power distribution**. Deformation-power and thrust-power distribution along the body of (a, b) swimmer $IS_d$, and (c, d) swimmer $SS_d$. The solid red line indicates the average over a single tail-beat period (from $t = 26$ to $t = 27$), whereas the envelope denotes the standard-deviation. The silhouettes at the bottom of each panel represent the fish body.

25

**Supplementary Movie S4.**   Detailed view of the flow-field around smart-swimmer $IS_\eta$. The top panel shows the vorticity field in colour and velocity vectors as black arrows. The middle panels show the swimming-efficiency and the deformation-power. The distribution of thrust-power and deformation-power along the swimmer's left-('lower') and right-lateral ('upper') surfaces are shown in the lower panels, and depict how these quantities depend on wake-interactions.

**Supplementary Movie S5.**   3D simulation of two nonautonomous swimmers, in which the leader swims steadily, and the follower maintains a specified relative position to interact favourably with the wake. The energetic-benefit for the follower is similar to that of each of the followers in Supplementary Movie S1.

**Supplementary Movie S6.**   3D simulation of three nonautonomous swimmers, in which the leaders use a feedback controller to maintain formation abreast of each other, and the follower holds a specified position relative to the leaders. The energetic-benefit for the follower is double that of the followers in Supplementary Movies 1 and 2, as it now interacts profitably with wake-rings generated by both the leaders.