



HAL
open science

Glider soaring via reinforcement learning in the field

Gautam Reddy, Jérôme Wong Ng, Antonio Celani, Terrence J Sejnowski,
Massimo Vergassola

► **To cite this version:**

Gautam Reddy, Jérôme Wong Ng, Antonio Celani, Terrence J Sejnowski, Massimo Vergassola. Glider soaring via reinforcement learning in the field. *Nature*, 2018, 562 (7726), pp.236-239. 10.1038/s41586-018-0533-0 . pasteur-02914599

HAL Id: pasteur-02914599

<https://pasteur.hal.science/pasteur-02914599>

Submitted on 12 Aug 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Soaring through reinforcement learning in the field

Gautam Reddy^{1*}, Jerome Wong Ng^{1*}, Antonio Celani², Terrence J. Sejnowski^{3,4}, Massimo Vergassola¹

¹Department of Physics, University of California San Diego, La Jolla, USA, ²The Abdus Salam International Center for Theoretical Physics, I-34014 Trieste, Italy, ³The Salk Institute for Biological Studies, La Jolla, USA and ⁴Division of Biological Sciences, University of California, San Diego, La Jolla, USA.

*These authors contributed equally to this work.

1 Abstract

2 Soaring birds often rely on ascending thermal plumes in the atmosphere as they search for
3 prey or migrate across large distances¹⁻⁴. The landscape of convective currents is rugged
4 and rapidly shifts on timescales of a few minutes as thermals constantly form, disintegrate,
5 or are transported away by the wind⁵⁻⁶. How soaring birds find and navigate thermals within
6 this complex landscape is unknown. Reinforcement learning⁷ provides an appropriate
7 framework to identify an effective navigational strategy as a sequence of decisions taken in
8 response to environmental cues. Here, we use reinforcement learning to train gliders in the
9 field to autonomously navigate atmospheric thermals. Gliders of two-meter wingspan were
10 equipped with a flight controller that enables an on-board implementation of autonomous
11 flight policies via precise control over their bank angle and pitch. A navigational strategy was
12 determined solely from the gliders' pooled experiences collected over several days in the
13 field using exploratory behavioral policies. The strategy relies on novel on-board methods to
14 accurately estimate the local vertical wind accelerations and the roll-wise torques on the
15 glider, which serve as navigational cues. We establish the validity of our learned flight policy
16 through field experiments, numerical simulations, and estimates of the noise in

17 measurements that is unavoidably present due to atmospheric turbulence. This is a novel
18 instance of learning a navigational task in the field, where learning is severely challenged by
19 a multitude of physical effects and the unpredictability of the natural environment. Our results
20 highlight the role of vertical wind accelerations and roll-wise torques as viable biological
21 mechanosensory cues for soaring birds, and provide a navigational strategy that is directly
22 applicable to the development of autonomous soaring vehicles.

23

24 Main

25 In reinforcement learning, an animal maximizes its long-term reward by taking actions in
26 response to its external environment and internal state. Learning occurs by reinforcing
27 behavior based on feedback from past experiences. Similar ideas have been used to
28 develop intelligent agents, reaching spectacular performance in strategic games like
29 backgammon⁸ and Go⁹, visual-based video game play¹⁰ and robotics^{11,12}. In the field,
30 physical constraints fundamentally prevent learning agents from using data-intensive
31 learning algorithms and the optimization of model design needed for quicker learning, which
32 are the conditions most often faced by living organisms.

33

34 A striking example in nature is provided by thermal soaring, where the extent of atmospheric
35 convection is not consistent across days and, even under suitable conditions, the locations,
36 sizes, durations and strengths of nearby thermals are unpredictable. As a result, the
37 statistics of training samples are skewed on any particular day. At smaller spatial and
38 temporal scales, fluctuations in wind velocities are due to turbulent eddies lasting a few
39 seconds that may mask or falsely enhance a glider's estimate of its mean climb rate.
40 Further, the measurement of navigational cues using standard instrumentation may be
41 consistently biased by aerodynamic effects, which requires precise quantification. Here, we
42 demonstrate that reinforcement learning can meet the challenge of learning to effectively
43 soar in atmospheric turbulent environments. To contrast with past work, the maneuvering of

44 an autonomous helicopter in ref. 11 is a control problem that is decoupled from
45 environmental fluctuations and has little trial-to-trial variability. Past autonomous soaring
46 algorithms have largely relied on locating the centroid of a drifting Gaussian thermal¹³⁻¹⁶,
47 which is unrealistic, or have applied learning methods in highly simplified simulated
48 settings¹⁷⁻¹⁹.

49

50 Using the reinforcement learning framework⁷, we may describe the behavior of the glider as
51 an agent traversing different states (s) by taking actions (a) while receiving a local reward (r).
52 The goal is to find a behavioral policy that maximizes the “value”, i.e., the mean sum of
53 future rewards up to a specified horizon. We seek a model-free approach, which estimates
54 the value of different actions at a particular state (called the Q function) solely through the
55 agent's experiences during repeated instances of the task, thereby bypassing the modeling
56 of complex atmospheric physics and aerodynamics (see Methods). The optimal policy is
57 subsequently derived by taking actions with the highest Q value at each state, where the
58 state includes sensorimotor cues and the glider's aerodynamic state.

59

60 To identify mechanosensory cues that could guide soaring, we recently combined above
61 ideas with simulations of virtual gliders in numerically generated turbulent flow²⁰. Two cues
62 emerged from our screening: (1) the vertical wind acceleration (\mathbf{a}_z) along the glider's path;
63 (2) the spatial gradients in the vertical wind velocity across the wings of the glider (ω).

64 Intuitively, the two cues correspond to the gradient of the vertical wind velocity in the
65 longitudinal and lateral directions of the glider, which locally orient it towards regions of
66 higher lift. Simulations in ref. 20 further showed that the glider's bank angle is the crucial
67 aerodynamic control variable; additional variables, such as the angle of attack, or other
68 mechanosensory cues, such as temperature or vertical velocity, offer minor improvements
69 when navigating within a thermal.

70

71 To learn to soar in the field, we used a glider (of two-meter wingspan) with autonomous
72 soaring capabilities (Figures 1A-B). The glider is equipped with a flight controller, which
73 implements a feedback control system used to modulate the glider's ailerons and elevator
74 such that a desired bank angle and pitch are maintained. Relevant measurements, such as
75 the altitude, ground velocity (\mathbf{u}), airspeed, bank angle (μ) and pitch, are made continuously
76 at 10 Hz using standard instrumentation (see Methods). At fixed time intervals, the glider
77 changes its heading by modulating its bank angle in accordance with the implemented
78 behavioral policy.

79

80 Noise and biases that affect learning in the field require the development of appropriate
81 methods to extract environmental cues from sensory devices' measurements. We found that
82 estimating \mathbf{a}_z by the derivative of the vertical ground velocity (\mathbf{u}_z), is significantly biased by
83 longitudinal motions of the glider about the pitch axis as the glider responds to an imbalance
84 of forces and moments while turning. By modeling the glider's longitudinal dynamics, we
85 obtain an unbiased estimate of the local vertical wind velocity (\mathbf{w}_z), and \mathbf{a}_z as its derivative
86 (Methods). The estimation of the spatial gradients across the wings, ω , poses a greater
87 challenge as it involves the difference between two noisy measurements at relatively close
88 positions. The key observation we used here is that the glider rolls due to contributions from
89 vertical wind velocity gradients, the feedback control mechanism and various aerodynamic
90 effects. The resulting roll-wise torque can be estimated from the small deviations of the true
91 bank angle from the desired one, and a novel dynamical model allows us to separate the ω
92 contribution due to velocity gradients from the other effects (Methods). A sample trace of the
93 resulting unbiased estimate of ω is shown in Figure 1C-D, together with traces of the vertical
94 wind velocity, \mathbf{w}_z , μ and unbiased estimates of \mathbf{a}_z .

95

96 Equipped with a proper procedure for estimating environmental cues, we next addressed the
97 specifics of learning in the field. First, to constrain our state space, we discretized the range
98 of values of \mathbf{a}_z and ω into three states each, positive high (+), neutral (0) and negative high (-
99). Second, we found that learning is accelerated by choosing \mathbf{a}_z attained at the subsequent
100 time step as the reward signal. The choice of \mathbf{a}_z (rather than \mathbf{w}_z) is an instance of reward
101 shaping that is justified in the Supplementary Information, where we show that using \mathbf{a}_z as a
102 reward still leads to a policy that optimizes the long-term gain in height. This property is a
103 special case of our general result that a particular reward function or its time derivatives (of
104 any order) yield the same optimal policy (Supplementary Information). Choosing \mathbf{w}_z as the
105 reward fails to drive learning in the soaring problem, possibly because the velocities (and
106 thus the rewards) are correlated across states and their temporal statistics strongly deviates
107 from the Markovianity assumption in reinforcement learning methods⁷. Indeed, velocity
108 fluctuations in turbulent flow are long-correlated, i.e. their correlation timescale is determined
109 by the largest timescale of the flow (see for instance Fig. 9 of ref. 21), which is of the order of
110 minutes in the atmosphere. Conversely, the correlation timescale of accelerations is
111 controlled by the smallest timescale²¹⁻²³ (the dissipation timescale in Fig. 7 of ref. 21). This is
112 estimated to be only a fraction of a second, which is much smaller than the time interval
113 between successive actions. Note that the previous experimental observations can be
114 rationalized by the combination of the power-law spectrum of turbulent velocity fluctuations
115 in the atmosphere and the extra factor of frequency squared in the spectrum of acceleration
116 vs velocity fluctuations²³. Finally, the glider's experiences, represented as state-action-state-
117 reward quadruplets, (s_t, a_t, s_{t+1}, r_t) , were cumulatively collected (over 15 days) into a set E
118 using explorative behavioral policies. Learning is monitored by bootstrapping the standard
119 deviation of the Q values from E (Figure 2A), calculated using value iteration methods
120 (Methods).

121

122 The navigational strategy derived at the end of the training period is presented in Figure 2B,
123 which shows the actions deemed optimal for the 45 possible states. Remarkably, the rows
124 corresponding to $\omega = 0$ resemble the so-called Reichmann rules²⁴ -- a set of simple heuristics
125 for soaring, which suggest a decrease/increase in bank angle when the climb rate
126 increases/decreases. Our strategy also gives a prescription for bank: for instance, when \mathbf{a}_z
127 and ω are both positive (top row in Figure 2B) i.e., in a situation when better lift is available
128 diagonal to the glider's heading, it is advantageous to bank not to the extreme but rather
129 maintain an intermediate value between -30° and -15° . Importantly, the learned
130 leftward/rightward bias in bank angle on encountering a positive/negative torque validates
131 our estimation procedure for ω .

132

133 In Figure 3A, we show a sample trajectory of the glider implementing the navigational
134 strategy in the field to remain aloft for ~ 12 minutes while spiraling to the height of low-lying
135 clouds (see also Extended Data Figure 1). On a day with strong atmospheric convection, the
136 time spent aloft is limited only by visibility and the receiver's range as the glider soars higher
137 or is constantly pushed away by the wind. A significant improvement in median climb rate of
138 0.35 m/s was measured in the field by performing repeated 3-minute trials over five days
139 (Figure 3B, Mann-Whitney $U = 429$, $n_{\text{control}} = 37$, $n_{\text{strategy}} = 49$, $p < 10^{-4}$ two-sided). Notably, this
140 value reflects a general improvement in performance averaged across widely variable
141 conditions without controlling for the availability of nearby thermals.

142

143 To examine possible advantages of larger gliders due to improved torque estimation, we
144 further analyzed soaring performance for different wingspans (l). While the naive expectation
145 is that the signal-to-noise ratio (SNR) in the estimation of ω scales linearly with l , we show
146 that the effects of atmospheric turbulence lead to a much weaker $l^{1/6}$ scaling (Methods).
147 Since testing our prediction would require a series of gliders with different wingspans, we

148 turned to numerical simulations of the convective boundary layer, adapted to reflect our
149 experimental setup (Methods). Results shown in Figure 3C-D are consistent with the
150 predicted scaling. Intuitively, the weak 1/6th exponent arises because the improvement in
151 gradient estimation is offset by the larger turbulent eddies, which only have a sweeping
152 effect for smaller wingspans, and contribute to velocity differences across the wings as l
153 increases. Our calculation yields an estimate of the SNR ~ 4 for typical experimental values;
154 similar arguments for \mathbf{a}_z yield an SNR ~ 7 . Experimental results, together with simulations
155 and SNR estimates, establish \mathbf{a}_z and ω as robust navigational cues for thermal soaring.

156

157 The real-world intricacies of soaring impose severe constraints on the complexity of the
158 underlying models, reflecting a fundamental trade-off between learning speed and
159 performance. Notably, the choice of a proper reward signal was crucial to make learning
160 feasible with the limited samples available. Though reward shaping has received some
161 attention in the machine learning community²⁵, its relevance for behaving animals remains
162 poorly understood. We remark that our navigational strategy constitutes a set of general
163 reactive rules with no learning performed during a particular thermal encounter. A soaring
164 bird may use a model-based approach of constantly updating its estimate of nearby
165 thermals' location based on recent experience and visual cues. Still, the importance of
166 vertical wind accelerations and torques for our policy suggests that they are likely useful for
167 *any* other strategy; our methods to estimate them in a glider suggest that they should be
168 accessible to birds as well. The hypothesis that birds utilize those mechanical cues while
169 soaring can be tested in experiments.

170

171 Finally, we note that single-thermal soaring is just one face of a multifaceted question: how
172 should a migrating bird or a cross-country glider fly among thermals over hundreds of
173 kilometers for a quick, yet risk-averse, journey²⁶⁻²⁸? This calls for the development of
174 effective methods for identifying areas of strong updraft based on mechanical and visual

175 cues. Such methods, coupled with our current work, pave the way towards a better
176 understanding of how birds migrate and the development of autonomous vehicles that can
177 extensively fly with minimal energy cost.

178

179 Main Figure Legends

180 **Figure 1: Soaring in the field using turbulent navigational cues.** (a) A trajectory of our
181 glider soaring in Poway, California. (b) A cartoon of the glider showing the available
182 navigational cues -- gradients in vertical wind velocities along the trajectory and across its
183 wings, which generate a vertical wind acceleration \mathbf{a}_z and a roll-wise torque ω respectively.
184 (c) A sample trace of the estimated vertical wind velocity \mathbf{w}_z and \mathbf{a}_z obtained in the field. (d)
185 The measured bank angle μ and the estimated ω during the same trial as in panel (c). The ω
186 (solid, green) is estimated from the small deviations of the measured bank angle (solid, blue)
187 from the expected bank angle (dashed, orange) after accounting for other effects (Methods).
188

189 **Figure 2: Convergence of the learning algorithm and the learned thermalling strategy.**
190 (a) The convergence of Q values during learning as measured by the standard deviation of
191 the mean Q value vs training time in the field, obtained by bootstrapping from the
192 experiences accumulated up to that point. (b) The final learned policy. Each symbol
193 corresponds to the best action (increasing/decreasing the bank angle μ by 15° or maintain
194 the same μ , as shown in the legend) to be taken when the glider observes a particular (\mathbf{a}_z, ω)
195 pair and is banked at μ . Combined symbols depict pairs of actions that are equally
196 rewarding. Note that a positive ω corresponds to a higher vertical wind velocity on the left
197 (right) wing of the glider and a positive (negative) μ corresponds to turning right (left) w.r.t the
198 glider's heading.

199

200 **Figure 3: Performance of the learned strategy and its dependence on the wingspan.**

201 (a) A 12-minute-long trajectory of the glider executing the learned thermalling strategy in the

202 field, colored by the vertical ground velocity u_z at each instant. (b) Experimentally measured
203 climb rate of a control random policy (black dots) is compared against the learned strategy
204 (red dots) over repeated 3-minute trials in the field. Each dot represents the average climb
205 rate in a single trial. A few outliers are not shown to restrict the range of the axis. (c)
206 Estimated signal-to-noise ratio (SNR) in ω and \mathbf{a}_z estimation vs wingspan (l) shown in green
207 and red respectively. The SNR for ω estimation is plotted in log-log scale (inset) to highlight
208 the weak $l^{1/6}$ scaling. (d) The mean climb rate for the learned strategy is compared for
209 different wingspans (red dots) in simulations of a glider soaring in the convective boundary
210 layer. For comparison we show the mean climb rates for a random policy and a strategy that
211 uses \mathbf{a}_z only (Methods). Error bars represent s.e.m.

212

213 Methods

214 **Experimental setup.** A Parkzone Radian Pro fixed-wing plane of 2-meter wingspan was
215 equipped with an on-board Pixfalcon autonomous flight controller operating on custom-
216 modified Arduplane firmware²⁹. The instrumentation available to the flight controller includes
217 a GPS, compass, barometer, airspeed sensor and an inertial measurement unit (IMU).
218 Measurements from multiple instruments are combined by an Extended Kalman Filter (EKF)
219 to give an estimate of relevant quantities such as the altitude z , the sink rate w.r.t ground
220 $-u_z$, pitch φ , bank angle μ and the airspeed V , at a rate of 10 Hz (see Extended Data Figure
221 2 for the definitions of the angles). Throughout the paper, we use $\mu > 0$ when the plane is
222 banked to the right and $\varphi > 0$ for the airplane pitched nose above the horizontal plane. For a
223 given desired pitch φ_d and desired bank angle μ_d , the controller modulates the aileron and
224 elevator control surfaces at 400 Hz using a proportional-integral-derivative (PID) feedback
225 control mechanism at a user-set time scale τ (see Extended Data Table 1 for parameter
226 values) such that:

227 $\tau \frac{d\phi}{dt} = \phi_d - \phi$ (1)

228 $\tau \frac{d\mu}{dt} = \mu_d - \mu$ (2)

229 ϕ_d is fixed during flight and can be used to indirectly modulate the angle of attack, α , which
 230 determines the airspeed and sink rate w.r.t air of the glider ($-\mathbf{v}_z$). Actions of increasing,
 231 decreasing or keeping the same bank angle are taken in time steps of t_a by changing the
 232 desired bank angle, μ_d , such that μ increases linearly from μ_i to μ_f in time interval t_a :

233 $\mu_d(t) = \mu_i + (\mu_f - \mu_i)(t + \tau)/t_a$ (3)

234

235 **Estimation of the vertical wind acceleration.** The vertical wind acceleration \mathbf{a}_z is defined
 236 as:

237 $\mathbf{a}_z \equiv \frac{d\mathbf{w}_z}{dt} = \frac{d}{dt}(\mathbf{u}_z - \mathbf{v}_z)$ (4)

238 where \mathbf{u} and \mathbf{v} are the velocities of the glider w.r.t the ground and air respectively, and \mathbf{w} is
 239 the wind velocity. Here, we have used the relation $\mathbf{w} = \mathbf{u} - \mathbf{v}$. An estimate of \mathbf{u} is obtained in
 240 a straightforward manner from the EKF, which combines the GPS and barometer readings to
 241 form the estimate. However, \mathbf{v}_z is confounded by various aerodynamic effects that
 242 significantly affect it on time scales of a few seconds (Extended Data Figure 3). Artificial
 243 accelerations introduced due to these effects impair accurate estimation of the wind
 244 acceleration and thus alter the perceived state during decision-making and learning. Two
 245 effects significantly influence variations in \mathbf{v}_z : (1) Sustained pitch oscillations with a period of
 246 a few seconds and varying amplitude, and (2) Angle of attack variations, which occur in
 247 order to compensate for the imbalance of lift and weight while rolling. In the Supplementary
 248 Information, we present a detailed analysis of the longitudinal motions that affect the glider,
 249 which is summarized here for conciseness. Changes in \mathbf{v}_z can be approximated as:

250 $\Delta v_z = -V(\Delta\alpha - \Delta\phi)$ (5)

251 where the Δ denotes the deviation from their value during steady, level flight. We obtain $\Delta\phi$
 252 directly from on-board measurements whereas $\Delta\alpha$ can be approximated for bank angle μ as:

253 $\Delta\alpha \approx (\alpha_0 - \alpha_i)\left(\frac{1}{\cos\mu} - 1\right)$ (6)

254 where α_0 is the angle of attack at steady, level flight and α_i is a parameter which depends on
 255 the geometry and the angle of incidence of the wing. The constant pre-factor ($\alpha_0 - \alpha_i$) is
 256 inferred from experiments. Measurements of \mathbf{u}_z together with the estimate of $\Delta\mathbf{v}_z$ are now
 257 used to estimate the vertical wind velocity \mathbf{w}_z up to a constant term, which can be ignored as
 258 it does not affect \mathbf{a}_z . The vertical wind acceleration \mathbf{a}_z is then obtained by taking the
 259 derivative of \mathbf{w}_z and is further smoothed using an exponential smoothing kernel of time scale
 260 σ_a (Extended Data Figure 4).

261

262 **Estimation of vertical wind velocity gradients across the wings.** Spatial gradients in the
 263 vertical wind velocity induce a roll-wise torque on the plane, which we estimate using the
 264 deviation of the measured bank angle from the expected bank angle. The total roll-wise
 265 torque on the plane has contributions from three sources – (1) the feedback control of the
 266 plane, (2) spatial gradients in the wind including turbulent fluctuations, and (3) roll-wise
 267 moments created due to various aerodynamic effects. Here, we follow an empirical
 268 approach: we note that the latter two contributions perturb the evolution of the bank angle
 269 from equation (2). We can then write an effective equation,

270 $\frac{d\mu}{dt} = \frac{\mu_d - \mu}{\tau} + \omega(t) + \omega_{aero}(t)$ (7)

271 where $\omega(t)$ and $\omega_{aero}(t)$ are contributions to the roll-wise angular velocity due to the wind and
 272 aerodynamic effects respectively. We empirically find four major contributions to ω_{aero} : (1)
 273 the dihedral effect, which is a stabilizing moment due to the effects of sideslip on a dihedral
 274 wing geometry, (2) the overbanking effect, which is a destabilizing moment that occurs
 275 during turns with small radii, (3) trim effects, which create a constant moment due to
 276 asymmetric lift on the two wings, and (4) a loss of rolling moment generated by the ailerons
 277 when rolling at low airspeeds. We quantify the contributions from the four effects and model
 278 their dependence on the bank angle (see Supplementary Information for more details on
 279 modeling and calibration). A estimate of ω is then obtained as:

280
$$\omega = \frac{d\mu}{dt} - \frac{\mu_d - \mu}{\tau} - \omega_{aero} \quad (8)$$

281 Finally, an exponential smoothing kernel is applied to obtain a smoothed ω (Extended Data
282 Figure 5).

283

284 **Design of the learning module.** The navigational component of the glider is modeled as a
285 Markov Decision Process (MDP), closely following the implementation used in ref. 20. The
286 Markovian transitions are discretized in time into intervals of size t_a . The state space consists
287 of the possible values taken by \mathbf{a}_z , ω and μ . To make the learning feasible within
288 experimental constraints and to maintain interpretability, we use a simple tile coding scheme
289 to discretize our state space: continuous values of \mathbf{a}_z and ω are each discretized into three
290 states (+,0,-), partitioned by thresholds $\pm K_a$, $\pm K_\omega$ respectively. The thresholds are set at \pm
291 0.8 times the standard deviation of \mathbf{a}_z and ω . Since the width of the distributions of \mathbf{a}_z and ω
292 can vary across days, the data obtained on a particular day is normalized by the standard
293 deviation calculated for that day. In effect, the filtration threshold to detect a signal against
294 turbulent “noise” is higher on days with more turbulence. The consequence is that the
295 behavior of the learned strategy could change across days, adapting to the recent statistics
296 of the environment. The bank angle takes five possible values -0° , $\pm 15^\circ$, $\pm 30^\circ$, while the
297 three possible actions allow for increasing, decreasing by 15° or keeping the same bank
298 angle. In summary, we have a total of $3 \times 3 \times 5 = 45$ states in the state space and 3 actions
299 in the action space.

300

301 We choose the local vertical wind acceleration \mathbf{a}_z obtained in the next time step as the
302 reward function. The choice of \mathbf{a}_z as an appropriate reward signal is motivated by
303 observations made in simulations from ref. 20. In the Supplementary Information, we show
304 that the obtained policy using \mathbf{a}_z as the reward function is equivalent to a policy that also
305 maximizes the expected gain in height.

306

307 **Learning the thermalling strategy in the field.** Data collected in the field is split into
308 (s,a,s',r) quadruplets containing the current state s , the current action a , the next state s' and
309 the obtained reward r , which are pooled together to obtain the transition matrix $T(s'|s,a)$ and
310 reward function $R(s,a)$. Value iteration methods are used to estimate the Q values from T
311 and R . The learning process is offline and off-policy; specifically, we begin training with a
312 'random' policy that takes the three possible actions with equal probability irrespective of the
313 current state as our behavioral policy, which was used for 12 out of the 15 days of training.
314 For the other days, a softmax policy⁷ with temperature set to 0.3 was used. For softmax
315 training, the Q values were first estimated from the data obtained in the previous days and
316 then normalized by the difference between the maximum and minimum Q values over the
317 three possible actions at a particular state, as described in ref. 20.

318

319 Using a fixed, random policy as our behavioral policy slows learning as state-action pairs
320 that rarely appear in the final policy are still sampled. On the other hand, calibrating the
321 parameters necessary for the unbiased measurement of \mathbf{a}_z and ω (see Supplementary
322 Information) is performed simultaneously with learning, which considerably reduces the
323 number of days required in the field. Importantly, offline learning permits us to continuously
324 monitor the variance of the estimated Q values by bootstrapping from the set E of
325 accumulated (s,a,s',r) quadruplets up to a particular point. Specifically, $|E|$ samples are
326 drawn with replacement from E and Q values are obtained for each state-action pair via
327 value iteration. The steps are repeated and the average of the bootstrapped standard
328 deviations in Q over all the state-action pairs is used as a measure of learning progress, as
329 shown in Figure 2A.

330

331 We expect certain symmetries in the transition matrix and the reward function, which we

332 exploit in order to expedite our learning process. Particularly, we note that the MDP is
333 invariant to an inversion of sign in the bank angle $\mu \rightarrow -\mu$. This transforms a state as $(\mathbf{a}_z, \omega, \mu)$
334 $\rightarrow (\mathbf{a}_z, -\omega, -\mu)$ and inverts the action from that of increasing the bank angle to decreasing the
335 bank angle and vice-versa. We symmetrize T and R as

$$336 \quad T^{sym} = \frac{T^+ + T^-}{2} \quad (9)$$

$$337 \quad R^{sym} = \frac{R^+ + R^-}{2} \quad (10)$$

338 where + and - denote the obtained values and those computed by applying the inverting
339 transformation respectively. Finally, T^{sym} and R^{sym} are used to obtain a symmetrized Q
340 function, which results in a symmetric policy as shown in Figure 2b. To conveniently obtain
341 the policy that uses only \mathbf{a}_z (Figure 3d), the above procedure is repeated with the threshold
342 for ω (K_ω) set to infinity.

343

344 **Testing the performance of the learned policy in the field.** To obtain the data shown in
345 Figure 3b, the glider is first sent autonomously to an arbitrary but fixed location 250 m above
346 ground level. The learned thermalling policy is then turned on and the mean climb rate i.e.,
347 the total height gained divided by the total time, is measured over a 3-minute interval. To
348 obtain the control data, the glider instead follows a random policy, which takes the three
349 possible actions with equal probability. The trials where we observe little to no atmospheric
350 convection were filtered out by imposing a threshold on the standard deviation of the vertical
351 wind velocity over the 3-minute trial. In Extended Data Figure 6, we show the distribution of
352 the standard deviation in w_z collected from ~240 3-minute trials over 9 days. Trials below the
353 threshold chosen as the 25th percentile mark (red, dashed line) are not used for our
354 analysis.

355

356 **Testing the performance for different wingspans in simulations.** Soaring performance is
357 analyzed in simulations similar to those developed in ref. 20 and adapted to reflect the
358 constraints faced by our glider and the environments typically observed in the field.

359

360 The atmospheric model consists of two components: (1) a kinematic model of turbulence
361 that reproduces the statistics of wind velocity fluctuations in the convective atmospheric
362 boundary layer, and (2) the positions, sizes and strengths of updrafts and downdrafts. The
363 temporal and spatial statistics of the generated velocity field satisfy the Kolmogorov and
364 Richardson laws³⁰ and the mean velocity profile in the convective boundary layer⁵, as
365 described in the SI of ref. 20. Stationary updrafts and downdrafts of Gaussian shape are
366 placed on a staggered lattice of spacing ~125m on top of the fluctuating velocity field.
367 Specifically, their contribution to the vertical wind velocity at position r is given by

368
$$w_z = \pm W e^{-\frac{(r_{\perp} - r_{\perp}^0)^2}{2R^2}} \quad (11)$$

369 where r_{\perp}^0 is the location of the center of the up(down)draft in the horizontal plane, W is its
370 strength and R is its radius. W is drawn from a half-normal distribution of scale 1.5m/s
371 whereas the radius is drawn from a (positive) normal distribution of mean 40m and deviation
372 10m. Gaussian white noise of magnitude ~0.2m/s is added as additional measurement
373 noise.

374

375 We assume the glider is in mechanical equilibrium; the lift, drag and weight forces on the
376 glider are balanced, except for centripetal forces while turning. The parameters
377 corresponding to the lift and drag curves and the (fixed) angle of attack are set such that the
378 airspeed is $V = 8\text{m/s}$ and the sink rate is 0.9m/s at zero bank angle, which match those
379 measured for our glider in the field. Control over bank angle is similar to those imposed in
380 the experiments i.e., the bank angle switches linearly between the angles 0° , $\pm 15^\circ$, $\pm 30^\circ$ in a
381 time interval t_a , corresponding to the time step between actions. The glider's trajectory and
382 wind velocity readings are updated every 0.1s. The vertical wind acceleration is derived

383 assuming that the glider directly reads the local vertical wind velocity. The vertical wind
384 velocity gradients across the wings are estimated as the difference between the vertical wind
385 velocities at the two ends of the wings. The readings are smoothed using exponential
386 smoothing kernels; the smoothing parameters in experiments are chosen to coincide with
387 those that yield the most gain in height in simulations.

388

389 **Estimation of the noise in gradient sensing due to atmospheric turbulence.** The cues
390 \mathbf{a}_z and ω measure the gradients in the vertical wind velocity along and perpendicular to the
391 heading of the glider. Updrafts and downdrafts are relatively stable structures in a varying
392 turbulent environment. Thermal detection through gradient sensing constitutes a
393 discrimination problem of deciding whether a thermal is present or absent given the current
394 \mathbf{a}_z and ω . We estimate the magnitude of turbulent ‘noise’ that unavoidably accompanies
395 gradient sensing. Intuitively, turbulent fluctuations in the atmospheric boundary layer (ABL)
396 are made up of eddies of different length scales, with the largest being the size of the height
397 of the ABL. Energy is transferred from larger, stronger eddies to smaller, weaker eddies, and
398 eventually dissipates at the centimeter scale due to viscosity in the bulk and the boundaries.
399 In the Supplementary Information, we present an explicit calculation of the signal to noise
400 ratio for ω estimation taking into account the effect of turbulent eddies on the statistics of
401 noise. Below, we give simple scaling arguments and refer to the Supplementary Information
402 for further details.

403

404 A glider moving at an airspeed V and integrating over a time scale T averages \mathbf{a}_z over a
405 length VT . For V much larger than the velocity scale of the eddies, which is typically the
406 case, the decorrelation of wind velocities is due to the glider’s motion; the eddies themselves
407 can be considered to be frozen in time. The magnitude of the spatial fluctuations across the
408 eddy of this size scales according to the Richardson-Kolmogorov law³⁰ as $\sim (VT)^{1/3}$. The
409 mean gradient signal when going up the gradient is $\sim (VT)$; the resultant signal to noise ratio
410 in \mathbf{a}_z scales as $(VT)^{2/3}$.

411

412 Similar arguments are applicable for ω measurements. In this case, the signal to noise ratio
413 has an additional dependence on the wingspan l . The dominant contribution to the noise
414 comes from eddies of size l , whose strength scales as $l^{1/3}$. As the glider moves a distance
415 VT , for $l \ll VT$, it traverses VT/l distinct eddies of size l . Consequently, the noise is averaged
416 out by a factor $(VT/l)^{-1/2}$, corresponding to the VT/l independent measurements. Multiplying
417 these two factors, the averaged noise is $\sim l^{5/6}(VT)^{-1/2}$. Since the mean gradient (i.e., the
418 signal) is $\sim l$, the signal to noise ratio is then $\sim l^{1/6}(VT)^{1/2}$.

419

420 From the above arguments and dimensional considerations, we get order-of-magnitude
421 estimates of the SNR for \mathbf{a}_z and ω estimation:

$$422 \quad SNR(\mathbf{a}_z) \sim \frac{WV^{2/3}T^{2/3}L^{1/3}}{wR} \quad (12)$$

$$423 \quad SNR(\omega) \sim \frac{WV^{1/2}T^{1/2}l^{1/6}L^{1/3}}{wR} \quad (13)$$

424 where W is the strength of the thermal, R is its radius, w is the magnitude of turbulent
425 vertical wind velocity fluctuations and L is the length scale of the ABL. For the SNR
426 estimates presented in the text, we use $W = 2\text{m/s}$, $R = 50\text{m}$, $l = 2\text{m}$, $V = 8\text{m/s}$, $T = 3\text{ s}$, $L = 1$
427 km . The values of V and T correspond to the airspeed of the glider in experiments and the
428 time scale between actions during learning respectively.

429

430 **Data availability.** The data that support the findings of this study are available from the
431 corresponding author upon reasonable request.

References

1. Newton I., *Migration Ecology of Soaring Birds*, Elsevier, 1st edition, 2008.
2. Shamoun-Baranes, J., Leshem, Y., Yom-tov, Y. and Liechti, O., Differential Use of Thermal Convection by Soaring Birds Over Central Israel, *The Condor*, 105(2): 208-218, 2003.
3. Weimerskirch, H., Bishop, C., Jeanniard-du-Dot, T., Prudor, A. and Sachs, G., Frigate birds track atmospheric conditions over months-long transoceanic flights, *Science*, 353(6294): 74-78, 2016.
4. Pennycuik, C. J., Thermal Soaring Compared in Three Dissimilar Tropical Bird Species, *Fregata Magnificens*, *Pelecanus Occidentals* and *Coragyps Atratus*, *J. Exp. Biol.*, 102:307-325, 1983.
5. Garrat, J. R., *The Atmospheric Boundary Layer*, Cambridge Atmospheric and Space Science Series, Cambridge University Press, 1994.
6. Lenschow, D. H., & Stephens, P. L., The role of thermals in the atmospheric boundary layer, *Boundary-Layer Meteorology*, 19:509-532, 1980.
7. Sutton, R. S., & Barto, A. G., *Reinforcement Learning: An Introduction*, MIT Press, 1st edition, 1998.
8. Tesauro, G., Temporal difference learning and TD-Gammon, *Commun ACM*, 38: 58-68, 1995.
9. Silver, D. et al, Mastering the game of Go without human knowledge, *Nature*, 550:354-359, 2017.
10. Mnih, V. et al, Human-level control through deep reinforcement learning, *Nature*, 518:529-533, 2014.
11. Kim, H. J., Jordan, M. I., Sastry, S., and Ng, A., Autonomous Helicopter Flight via Reinforcement Learning, *Advances in Neural Information Processing Systems*, 16, 2003.

12. Levine, S., Finn, C., Darrell, T. and Abbeel, P., End-to-End Training of Deep Visuomotor Policies, *Journal of Machine Learning Research*, 17:1-40, 2016.
13. Allen, M. J. & Lin, V., Guidance and Control of an Autonomous Soaring UAV, *Proceedings of the AIAA Aerospace Sciences Meeting*, (American Institute of Aeronautics and Astronautics, Reston, VA), AIAA Paper 2007-867.
14. Edwards, D. J., Implementation Details and Flight Test Results of an Autonomous Soaring Controller, AIAA Guidance, Navigation and Control Conference and Exhibit, 2008.
15. Edwards, D. J., Autonomous Soaring: The Montague Cross Country Challenge
Doctorate theses, North Carolina State University, Aerospace Engineering, Raleigh, North Carolina, 2010.
16. Ákos, Z., Nagy, M., Leven, S., and Vicsek, T., Thermal soaring flight of birds and unmanned aerial vehicles. *Bioinspiration & Biomimetics*, 5(4), 045003, 2010.
17. Doncieux, S., Mouret, J. B. and Meyer J-A., Soaring behaviors in UAVs : 'animat' design methodology and current results. In *7th European Micro Air Vehicle Conference (MAV07)*, Toulouse, 2007. □
18. Wharington, J. and Herszberg, I., Control of a high endurance unmanned aerial vehicle, *Proceedings of the 21st Congress of International Council of the Aeronautical Sciences* (International Council of the Aeronautical Sciences, Bonn, Germany), Paper 98-3.7.1, 1998. □
19. Chung J. J., Lawrance, N. R. J. and Sukkarieh, S, Learning to soar: Resource-constrained exploration in reinforcement learning, *The International Journal of Robotics Research*, 34(2):158-172, 2014.
20. Reddy, G., Celani, A., Sejnowski, T. and Vergassola, M., Learning to soar in turbulent environments, *Proc. Natl. Acad. Sci.*, 113(33): E4877-E4884, 2016.
21. Yeung, P. K. and Pope, S. B., Lagrangian statistics from direct numerical simulations of isotropic turbulence, *J. Fluid. Mech.*, 207:531-586, 1989.

22. Voth, G. A., La Porta, A., Crawford, A. M., Alexander, J., and Bodenschatz, E., Measurement of particle accelerations in fully developed turbulence, *J. Fluid. Mech.*, 469:121-160, 2002.
23. Tennekes, H. and Lumley, J. L., *A first course in turbulence*, MIT Press, 1972.
24. Reichmann, H., *Cross-Country Soaring*, Thomson Publications, Santa Monica, CA, 1988.
25. Ng, A. Y., Harada, D., and Russell, S. J., Policy Invariance Under Reward Transformations: Theory and Application to Reward Shaping, *Proc. of the 16th International Conference on Machine Learning*, P278-287, 1999.
26. MacCready, P. B. J., Optimum airspeed selector, *Soaring*, 1958(1): 10–11, 1958.
27. Horvitz, N. et al., The gliding speed of migrating birds: Slow and safe or fast and risky? *Ecology Letters*, 17(6):670–679, 2014.
28. Cochrane, J. H., MacCready theory with uncertain lift and limited altitude, *Technical Soaring*, 23:88–96, 1999.
29. ArduPilot, www.ardupilot.org, (2018).
30. Frisch, U., *Turbulence: The Legacy of A. N. Kolmogorov*, Cambridge University Press, 1995.

Supplementary Information

Supplementary Information is linked to the online version of the paper at

www.nature.com/nature.

Acknowledgements

This work was supported by Simons Foundation Grant 340106 (to M.V.).

Author Contributions

All authors were involved in designing the study and drafting the final manuscript. G.R. and J.W.N. performed the experiments and analyzed the data. G.R., A.C. and M.V. contributed to the theoretical results.

Author Information

Reprints and permissions information is available at www.nature.com/reprints.

The authors declare no competing financial interests.

Correspondence and requests for materials should be addressed to M.V.
(massimo@physics.ucsd.edu).

Extended Data Legends

Extended Data Table 1: Parameter values.

Extended Data Figure 1: Sample trajectories obtained in the field. The three-dimensional view and top view of the glider's trajectory as it executes the learned thermalling strategy (labeled 's') or a random policy that takes actions with equal probability (labeled 'r'). The trajectories are colored with the instantaneous vertical ground velocity (u_z). The green (red) dot shows the start (end) point of the trajectory. Trajectories s1, s2 and r1 last for 3 minutes each, whereas s3 lasts for ~8 minutes.

Extended Data Figure 2: Force-body diagram of a glider. The forces on a glider and the definitions of the various angles that determine the glider's motion.

Extended Data Figure 3: Modeling the longitudinal motion of the glider. (a) A sample trajectory of a glider's pitch and its vertical velocity w.r.t ground u_z in a case where the feedback control over the pitch is reduced in order to exaggerate the pitch oscillations. The blue line shows the measured u_z and the orange line is u_z obtained after subtracting the contributions from longitudinal motions of the glider (see Supplementary Information). (b) The blue line shows the average change in u_z when a particular action is taken (labeled above each panel), averaged over n three-second intervals. The 13 panels correspond to the 13 possible bank angle changes from the angles 0° , $\pm 15^\circ$, and $\pm 30^\circ$ by increasing, decreasing the bank angle by 15° or keeping the same angle. The green, dashed line shows the prediction from the model whereas the orange line is the estimated w_z . The axis on the right shows the averaged pitch as a red, dashed line.

Extended Data Figure 4: The estimated vertical wind acceleration is unbiased after accounting for the glider's longitudinal motion. (a) The averaged vertical wind

acceleration, \mathbf{a}_z in units of its standard deviation \mathbf{a}_z , plotted as in Extended Data Figure 3b, is shown in orange with (blue line) and without (orange line) accounting for the glider's longitudinal motions. The axis on the right shows the airspeed as a green, dashed line. (b) The PDFs (probability density functions) of \mathbf{a}_z for the different bank angle changes. The black, dashed line shows the median.

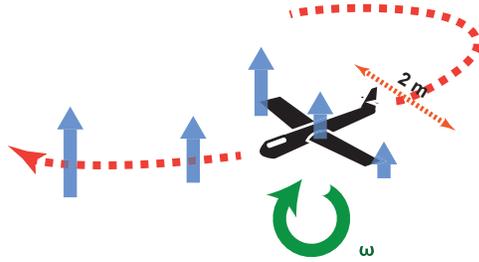
Extended Data Figure 5: The estimated roll-wise torque is unbiased after accounting for the effects of feedback control and glider aerodynamics. (a) The averaged evolution of the bank angle shown as in Extended Data Figure 3b. The blue line shows the measured bank angle and the dashed, orange line shows the best-fit line obtained from simultaneously fitting the 13 blue curves to the prediction (see Supplementary Information). (b) The PDFs (probability density functions) of the roll-wise torque ω (in units of its standard deviation) for the different bank angle changes. The black, dashed line shows the median value.

Extended Data Figure 6: The distribution of the strength of vertical currents observed in the field. The root-mean-square vertical wind velocity measured in the field is pooled from ~240 3-minute trials collected over 9 days. The dashed, red line shows the threshold criterion imposed when measuring the performance of the strategy in the field (see Methods).

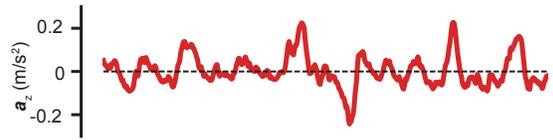
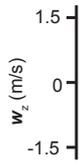
a



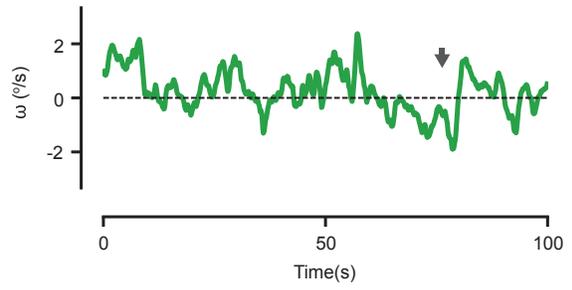
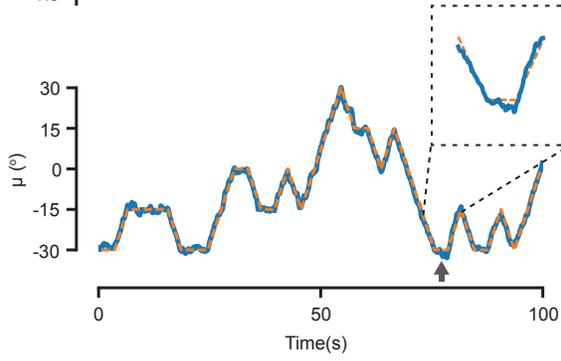
b

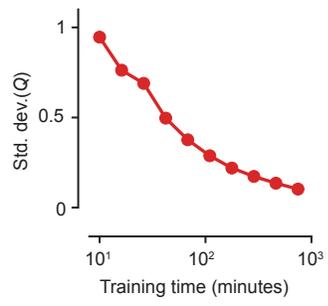


c



d



a**b**