

Q-Learning to navigate turbulence without a map.

Marco Rando^{a1}, Martin James^b, Alessandro Verri^a, Lorenzo Rosasco^a, Agnese Seminara^{b1}

April 29, 2024

Abstract

We consider the problem of olfactory searches in a turbulent environment. We focus on agents that respond solely to odor stimuli, with no access to spatial perception nor prior information about the odor location. We ask whether navigation strategies to a target can be learned robustly within a sequential decision making framework. We develop a reinforcement learning algorithm using a small set of interpretable olfactory states and train it with realistic turbulent odor cues. By introducing a temporal memory, we demonstrate that two salient features of odor traces, discretized in few olfactory states, are sufficient to learn navigation in a realistic odor plume. Performance is dictated by the sparse nature of turbulent plumes. An optimal memory exists which ignores blanks within the plume and activates a recovery strategy outside the plume. We obtain the best performance by letting agents learn their recovery strategy and show that it is mostly casting cross wind, similar to behavior observed in flying insects. The optimal strategy is robust to substantial changes in the odor plumes, suggesting minor parameter tuning may be sufficient to adapt to different environments.

Bacterial cells localize the source of an attractive chemical even if they hold no spatial perception. They respond solely to temporal changes in chemical concentration and the result of their response is that they move toward attractive stimuli by climbing concentration gradients [1]. Larger organisms also routinely sense chemicals in their environment to localize or escape targets,

but cannot follow chemical gradients since turbulence breaks odors into sparse pockets and gradients lose information [2, 3, 4, 5, 6]. The question of which features of turbulent odor traces are used by animals for navigation is natural, but not well understood. Beyond olfaction, some animals could use also prior spatial information to navigate [7, 8, 9, 10], but if and how chemosensation and spatial perception are coupled is still not clear.

An algorithmic perspective to olfactory navigation in turbulence can shed light on some of these questions. Without aiming at an exhaustive taxonomy, see e.g. [11] for a recent review, we recall some approaches relevant to put our contribution in context. One class of methods are biomimetic algorithms, where explicit navigation rules are crafted taking inspiration from animal behavior. An advantage of these methods is interpretability, in the sense that they provide insights into features that effectively achieve turbulent navigation, for example: odor presence/absence [12, 13, 14, 5]; odor slope at onset of detection [15]; number of detections in a given interval of time [16] and the time of odor onset [17]. On the flip side, in biomimetic algorithms behaviors are hardwired and typically reactive, not relying on any optimality criterion.

A way to tackle this shortcoming is to cast olfactory navigation within a sequential decision making framework [18]. In this context, navigation is formalized as a task with a reward for success; by maximizing reward, optimal strategies can be sought to efficiently reach the target. A byproduct is that most algorithmic choices can often be done in a principled way. Within this framework, some approaches make explicit use of spatial information. Bayesian algorithms use a spatial map to guess the target location and use odor to refine this guess or “belief”. A prominent algorithm for olfactory navigation based on the concept of belief is the information-seeking algorithm [3] akin to exploration heuristics widely used in robotics [19, 20] (see e.g. [21, 22]). Using Bayesian sequential decision making and the notion of beliefs, navigation can be formalized as a Partially Observable Markov Decision Process (POMDP) [23, 24, 25], that can be approximatively solved [26, 27, 28]. POMDP

^aMalGa - DIBRIS, University of Genova, Italy

^bMalGa - DICCA, University of Genova, Italy

The authors declare no conflict of interest.

¹To whom correspondence should be addressed. E-mail: marco.rando@edu.unige.it, agnese.seminara@unige.it

Author’s contributions. M.R., A.V., L.R. and A.S. designed research; M.R., M.J., A.V., L.R. and A.S. performed research; M.R., L.R. and A.S. wrote manuscript.

Keywords: Navigation | Reinforcement learning | Olfaction | Turbulence | Memory

approaches are appealing since beliefs are a sufficient statistics for the entire history of odor detections. However, they are computationally cumbersome. Further, they leave the question open of whether navigation as sequential decision making can be performed using solely olfactory information.

Recently, two algorithms studied navigation as a response to olfactory input alone [29, 30]. In [29] artificial neural networks were shown to learn near optimal strategies, but they were trained on odor cues with limited sparsity, and training with sparse odor cues typical of turbulence remains to be tested. In [30] an approach based on finite state controllers was proposed. Here, optimization was done using a model-based technique, relying on prior knowledge of the likelihood to detect the odor in space, hence still using spatial information. A different model free optimization could also be considered avoiding spatial information but this latter approach also remains to be tested. More generally, all the above approaches manipulate internally the previous history (memory) of odor detections. In this sense they are less interpretable, since the features of odor traces that drive navigation do not emerge explicitly.

In this paper, we propose a reinforcement learning (RL) approach to navigation in turbulence based on a set of interpretable olfactory features, with no spatial information, and highlight the role played by memory within this context. More precisely, we learn optimal strategies from data by training tabular Q learning [18] with realistic odor cues obtained from state-of-the-art Direct Numerical Simulations of turbulence. From the odor cues, we define features as moving averages of odor intensity and sparsity: the moving window is the temporal memory and naturally connects to the physics of turbulent odors. States are then obtained discretizing such features. Due to sparsity, agents may detect no odor within the moving window. We show there is an optimal memory minimizing the occurrence of this “void state”. The optimal memory scales with the blank time dictated by turbulence as it emerges from a trade off requiring that: *(i)* short blanks— typical of turbulent plumes— are ignored by responding to detections further in the past, and *(ii)* long blanks promptly trigger a recovery strategy to make contact with the plume again. We leverage these observations to tune the memory adaptively, by setting it equal to the previous blank experienced along an agent’s path. With this choice, the algorithm tests successfully in distinct environments, suggesting that tuning can be made robustly to enable generalization. The agent learns to surge upwind in most non-void states and to recover by casting crosswind in the absence of detections. Optimal agents limit encounters with the void state to a narrow band right at the edge of the plume. This suggests

that the temporal odor features we considered effectively predict when the agent is exiting the plume and point to an intimate connection between temporal predictions and spatial navigation.

Significance

Finding mates or food in the presence of turbulence is challenging because odors constantly switch on and off unpredictably. As a result, it is unclear whether animals couple odor to other sources of information, what are the relevant features of odor stimuli and how they change according to the environment. A long history of bioinspired algorithms address this problem by crafting rules for navigation that mimic animal behavior: but can effective navigation be learned from the environment rather than set a priori? To address this question we train a reinforcement learning algorithm with realistic turbulent stimuli. Searchers learn to navigate by trial and error and respond solely to odor, with no further input. All computations are defined explicitly, enhancing interpretability. The upshot is that the algorithm identifies odor features as averages over a temporal scale (memory) dictated by the time between odor detections and thus by physics. There is no need to know physics beforehand, as memory can be adjusted based on experience. This approach naturally complements previous algorithms that use prior information and a map of space to plan navigation, rather than learn it from the environment. To what extent different animals plan *vs* learn to navigate remains to be understood.

Results

Background Given a source of odor placed in an unknown position of a two-dimensional space, we consider the problem of learning to reach the source, Figure 1A. We formulate the problem as a discrete Markov Decision Process by discretizing space in tiles, also called “grid-world” in the reinforcement learning literature [18]. In this problem, an agent is in a given state s which is one of a discrete set of n tiles: $s \in S := \{s_1, \dots, s_n\}$. At each time step it chooses an action a which is a step in any of the coordinate directions $a \in \{\text{up,down,left,right}\}$. The goal is to find sequences of actions that lead to the source as fast as possible and is formalized with the notion of policy and reward which we will introduce later. If agents have perfect knowledge of their own location and of the location of the source in space, the problem reduces to finding the shortest path.

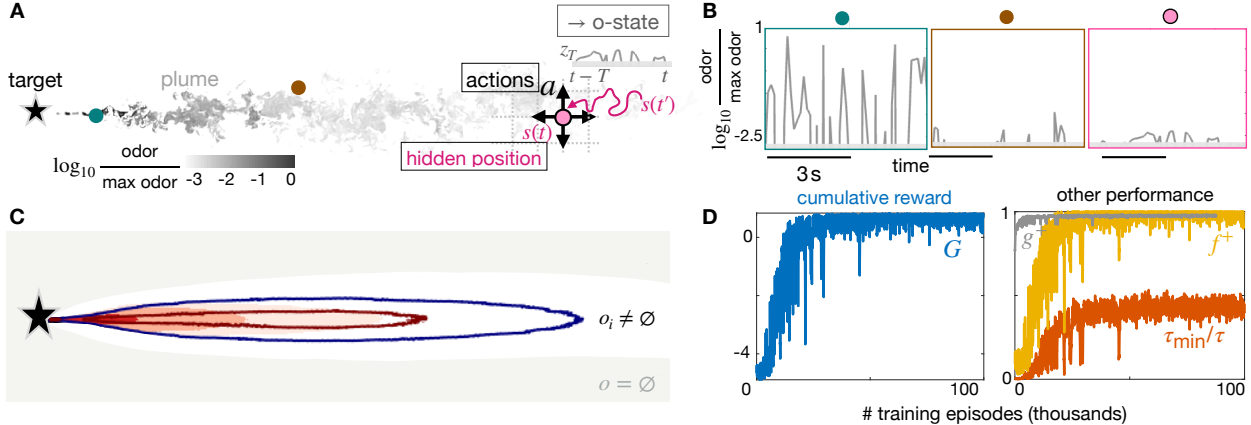


Figure 1: Learning a stimulus-response strategy for turbulent navigation. (A) Representation of the search problem with turbulent odor cues obtained from Direct Numerical Simulations of fluid turbulence (grey scale, odor snapshot from the simulations). The discrete position s is hidden; the odor concentration $z_T = z(s(t'), t') | t - T \leq t' \leq t$ is observed along the trajectory $s(t')$, where T is the sensing memory. (B) Odor traces from direct numerical simulations at different (fixed) points within the plume. Odor is noisy and sparse, information about the source is hidden in the temporal dynamics. (C) Contour maps of olfactory states with nearly infinite memory ($T = 2598$): on average olfactory states map to different locations within the plume and the void state is outside the plume. (D) Performance of stimulus-response strategies obtained during training, averaged over 500 episodes. We train using realistic turbulent data with memory $T = 20$ and backtracking recovery.

Using time *vs* space to address partial observability. In our problem however, the agent does not know where the source is, hence its position s relative to the source, is unknown or “partially observed”. Instead, it can sense odor released by the target. In the language of RL, odor is an “observation” – but does it hold information about the position s ? The answer is yes: several properties of odor stimuli depend on the distance from the source. However in the presence of turbulence, information lies in the statistics of the odor stimulus. Indeed, when odor is carried by a turbulent flow, it develops into a dramatically stochastic plume, i.e. a complex and convoluted region of space where the fluid is rich in odor molecules. Turbulent plumes break into structures that distort and expand while they travel away from their source and become more and more diluted [31, 4, 32, 6], see Figure 1A. As a consequence, an agent within the plume experiences intermittent odor traces that endlessly switch on (whiff) and off (blank) Figure 1B. The intensity of odor whiffs and how they are interleaved with blanks depends on distance from release, as dictated by physics [32]. Thus the upshot of turbulent transport is that the statistical properties of odor traces depend intricately on the position of the agent relative to the source. In other words, information about the state s is hidden within the observed odor traces.

This positional information can be leveraged with a Bayesian approach that relies on guessing s , i.e. defining the probability distribution of the position, also called belief. This is the approach that has been more commonly adopted in the literature until now [26, 27, 28]. Note that because of the complexity of these algorithms, only relatively simple measures of the odor are computationally feasible, e.g. instantaneous presence/absence. Here we take a different model-free and map-free approach. Instead of guessing the current state s , we ignore the spatial position and respond directly to the temporal traces of the odor cues. Two other algorithms have been proposed to solve partial observability by responding solely to odor traces with recurrent neural networks [29] and finite state controllers [30] that manipulate implicitly the odor traces. Here instead we manipulate odor traces explicitly, by defining memory as a moving window and by crafting a small number of features of odor traces.

Features of odor cues: definition of discrete olfactory states and sensing memory. To learn a response to odor traces, we first define a finite set of *olfactory states*, $o \in O$, so that they bear information about the location s . Defining the olfactory states is a challenge due to the dramatic fluctuations and irregularity of turbulent odor traces. To construct a fully interpretable

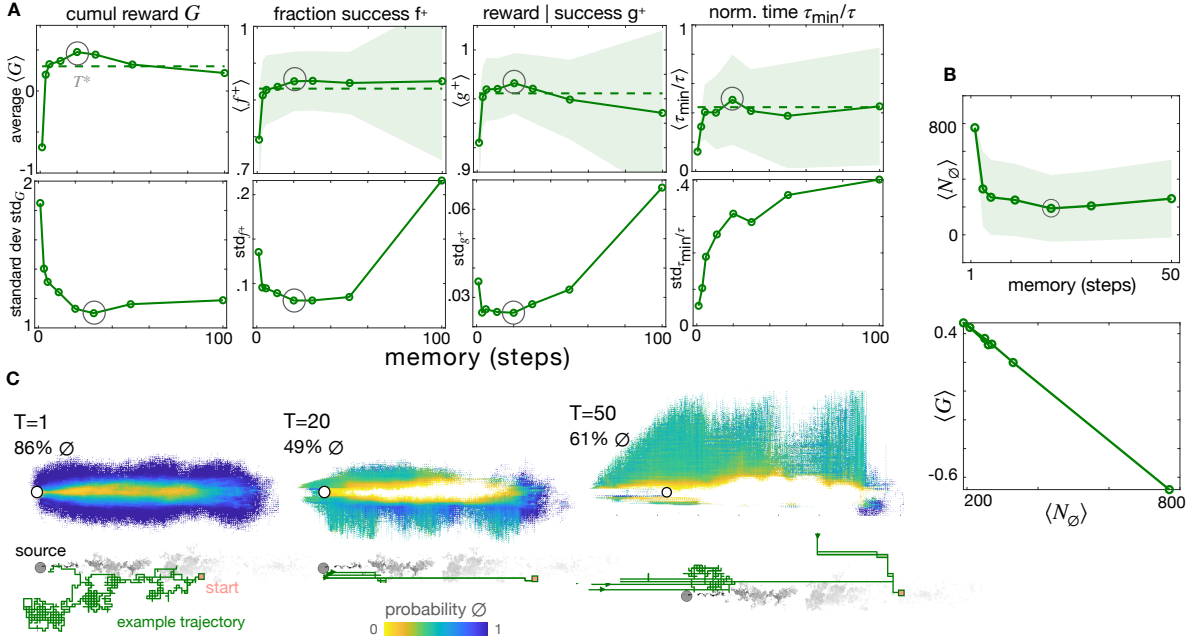


Figure 2: The optimal memory T^* . (A) Four measures of performance as a function of memory with backtracking recovery (solid line) show that the optimal memory $T^* = 20$ maximizes average performance and minimizes standard deviation, except for the normalized time. Top: Averages computed over 10 realizations of test trajectories starting from 43000 initial positions (dash: results with adaptive memory). Bottom: standard deviation of the mean performance metrics for each initial condition (see Materials and Methods). (B) Average number of times agents encounter the void state along their path, $\langle N_\emptyset \rangle$, as a function of memory (top); cumulative average reward $\langle G \rangle$ is inversely correlated to $\langle N_\emptyset \rangle$ (bottom), hence the optimal memory minimizes encounters with the void. (C) Colormaps: Probability that agents at different spatial locations are in the void state at any point in time, starting the search from anywhere in the plume and representative trajectory of a successful searcher (green solid line) with memory $T = 1$, $T = 20$, $T = 50$ (left to right). At the optimal memory agents in the void state are concentrated near the edge of the plume. Agents with shorter memories encounter voids throughout the plume; agents with longer memories encounter more voids outside of the plume as they delay recovery. In all panels, shades are \pm standard deviation.

low dimensional state space, we aim at a small number of olfactory states that bear robust information about s , i.e. for all values of s . We previously found that pairing features of sparsity as well as intensity of turbulent odor traces predicts robustly the location of the source for all s [33]. Guided by these results, we use these two features extracted from the temporal history of odor detections to define a small set of olfactory states.

We proceed to define a function that takes as input the history of odor detections along an agent's path and returns its current olfactory state. We indicate with $s(t)$ the (unknown) path of an agent, and with $z(s(t), t)$ the observations i.e. odor detections along its path. First, we define a sensing memory T and we consider a short excerpt of the history of odor detections z_T of duration T prior to the current time t . Formally,

$z_T(t) := \{z(s(t'), t') \mid t - T \leq t' \leq t\}$. Second, we measure the average intensity c (moving average of odor intensity over the time window T , conditioned to times when odor is above threshold), and intermittency i (the fraction of time the odor is above threshold during the sensing window T). Both features c and i are described by continuous, positive real numbers. Third, we define 15 olfactory states by discretizing i and c in 3 and 5 bins respectively. Intermittency i is bounded between 0 and 1 and we discretize it in 3 bins by defining two thresholds (33% and 66%). The average concentration, c , is bounded between 0 and the odor concentration at the source, hence prior information on the source is needed to discretize c using set thresholds. To avoid relying on prior information, we define thresholds of intensity as percentiles, based on a histogram that is populated on-

line, along each agent’s path (see Materials and Methods). The special case where no odor is detected over T deserves attention, hence we include it as an additional state named “void state” and indicate it with $o \equiv \emptyset$. When T is sufficiently long, the resulting olfactory states map to different spatial locations (Figure 1C, with T equal to the simulation time). Hence this definition of olfactory states can potentially mitigate the problem of partial observability using temporal traces, rather than spatial maps. But will these olfactory states with finite memory T guide agents to the source?

Q learning: a map-less and model-free navigation to odor sources. To answer this question, we trained tabular episodic Q learning [18]. In a nutshell, we use a simulator to place an agent at a random location in space at the beginning of each episode. The agent is not aware of its location in space, but it senses odor provided by the fluid dynamics simulator and thus can compute its olfactory state o , based on odor detected along its path in the previous T sensing window. It then makes a move according to a set policy of actions $a \sim \pi_0(a|o)$. After the move, the simulator displaces the agent to its new location and relays the agent a negative reward $R = -\eta$ if it is not at the source and a positive reward $R = 1$ if it reaches the source. The goal of RL is to find a policy of actions that maximizes the expected cumulative future reward $G = E_\pi(\sum_{t=0}^{\infty} \gamma^t R_{t+1})$ where the expectation is over the ensemble of trajectories and rewards generated by the policy from any initial condition. Because reward is only positive at the source, the optimal policy is the one that reaches the source as fast as possible. To further encourage the agent to reach the source quickly, we introduce a discount factor $\gamma < 1$.

Episodes where the agent does not reach the source are ended after $H_{\max} = 5000$ with no positive reward. As it tries actions and receives rewards, the agent learns how good the actions are. This is accomplished by estimating the quality matrix $Q(o, a)$, i.e. the maximum expected cumulative reward conditioned to being in o and choosing action a at the present time: $Q(o, a) = \max_\pi E_\pi(\sum_{t=0}^{\infty} \gamma^t R_{t+1} | o_t = o, a_t = a)$. At each step, the agent improves its policy by choosing more frequently putatively good actions. Once the agent has a good approximation of the quality matrix, the optimal policy corresponds to the simple readout: $\pi^*(a|o) = \delta(a - a^*(o))$ where $a^*(o) = \arg \max_a Q(o, a)$, for non-void states $o \neq \emptyset$.

Recovery strategy. To fully describe the behavior of our Q-learning agents, we have to prescribe their policy from the void state $o \equiv \emptyset$. This is problematic

because turbulent plumes are full of holes thus the void state can occur anywhere both within and outside the plume, Figure 1A. As a consequence, the optimal action $a^*(\emptyset)$ from the void state is ill defined. We address this issue by using a separate policy called “recovery strategy”. Inspired by path integration as defined in biology [34, 35, 36], we propose the backtracking strategy consisting in retracing the last T_a steps after the agent lost track of the odor. If at the end of backtracking the agent is still in the void state, it activates Brownian motion. Backtracking requires that we introduce memory of the past T_a actions. This timescale T_a for activating recovery is conceptually distinct from the duration of the sensing memory – however here we set $T_a = T$ for simplicity.

We find that Q-learning agents successfully learn to navigate to the odor source by responding solely to their olfactory state, with no sense of space nor models of the odor cues. Learning can be quantified by monitoring the cumulative reward which continuously improves with further training episodes (Figure 1D, left). Improved reward corresponds to agents learning how to reach the source more quickly and reliably with training. Indeed, it is easy to show that the expected cumulative reward $G = \langle e^{-\lambda\tau} - \eta(1 - e^{-\lambda\tau})/(1 - \gamma) \rangle$, where τ is a random variable corresponding to time to reach the source and $\gamma = e^{-\lambda\Delta t}$ is the discount factor, with the time step $\Delta t = 1$ (see Materials and Methods). Large rewards arise when (i) a large fraction f^+ of agents successfully reaches the source and (ii) the agents reach the source quickly, which maximizes $g^+ = \langle e^{-\lambda\tau} | \text{success} \rangle$. Indeed $G = f^+G^+ + (1 - f^+)G^-$, where $G^+ = g^+ - \eta(1 - g^+)/(1 - \gamma)$ and $G^- = -\eta(1 - e^{-\lambda H_{\max}})/(1 - \gamma)$. H_{\max} is the horizon of the agent i.e. the maximum time the agent is allowed to search, and after which the search is considered failed. Note that agents starting closer to the target receive larger rewards purely because of their initial position. To monitor performance independently on the starting location, we introduce the inverse time to reach the source relative to the shortest-path time from the same initial location, which goes from 0 for failing agents to 1 for ideal agents $\langle \tau_{\min}/\tau \rangle$, independently on their starting location. Note that this is not the quantity that is optimized for. One may specifically target this performance metrics, which is agnostic on the duration of an agent’s path, by discounting proportionally to t/τ_{\min} . All four measures of performance plateau to a maximum, suggesting learning has achieved a nearly optimal policy (Figure 1D). Once training is completed, we simulate the trajectory of test-agents starting from any of the about 43 000 admissible locations within the plume and moving according to the optimal policy. We will reca-

pitulate performance with the cumulative reward G averaged over the test-agents and dissect it into speed g^+ , convergence f^+ and relative time $\langle \tau_{\min}/\tau \rangle$.

Optimal memory. By repeating training using different values of T we find that performance depends on memory and an optimal memory T^* exists (Figure 2A). Why is there an optimal memory? The shortest memory $T = 1$ corresponds to instantaneous olfactory states: the instantaneous contour maps of the olfactory states are convoluted and the void state is pervasive (Figure 2C, top). As a consequence, agents often activate recovery even when they are within the plume, the policy almost always leads to the source ($f^+ = 79\% \pm 13\%$) but follows lengthy convoluted paths ($\tau_{\min}/\tau = 0.14 \pm 0.05$, Figure 2A, bottom). As memory increases, the olfactory states become smoother and agents encounter less voids (Figure 2C, center), perform straighter trajectories ($\tau_{\min}/\tau = 0.5 \pm 0.3$) and reach the source reliably ($f^+ = 95\% \pm 8\%$), Figure 2C, bottom. Further increasing memory leads to even less voids within the plume and even smoother olfactory states. However – perhaps surprisingly – performance does not further improve but slightly decreases (at $T = 50$, $f^+ = 94\% \pm 8\%$ and $\tau_{\min}/\tau = 0.38 \pm 0.36$). A long memory is deleterious because it delays recovery from accidentally exiting the plume, thus increases the number of voids *outside* of the plume (Figure 2C, bottom). Indeed, agents often leave the plume accidentally as they measure their olfactory state *while they move*. They receive no warning, but realize their mistake after T steps, when they enter the void state and activate recovery to re-enter the plume. The delay is linear with memory when agents recover by backtracking, but it depends on the recovery strategy (see discussion below and Supplementary Figure 1).

Thus short memories increase time in void *within* the plume, whereas long memories increase time in void *outside* the plume: the optimal memory minimizes the overall chances to experience the void (Figure 2B). Intuitively, T^* should match the typical duration τ_b of blanks encountered within the plume, so that voids within the plume are effectively ignored without delaying recovery unnecessarily. Consistently, $\langle \tau_b \rangle$ averaged across all locations and times within the plume is $\langle \tau_b \rangle = 9.97 \pm 41.16$, comparable with the optimal memory T^* (Figure 2A).

Adaptive memory. There is no way to select the optimal memory T^* without comparing several agents or relying on prior information on the blank durations. In order to avoid prior information, we venture to de-

fine memory adaptively along each agent’s path, using the intuition outlined above. We define a buffer memory T_b , and let the agent respond to a sensing window $T < T_b$. Ideally we would like to set $T \sim \langle \tau_b \rangle$. With this choice, blanks shorter than the average blank are ignored, as they are expected within the plume, whereas blanks longer than average initiate recovery, as they signal that the agent exited the plume. However, agents do not have access to $\langle \tau_b \rangle$ hence we set $T = \tau_b^-$, where τ_b^- is the most recent blank experienced by the agent. With this choice, the sensing memory T fluctuates considerably along an agent’s path, due to turbulence ([32] and Figure 3A-B). Note that blanks are estimated along paths, thus the statistics of T only qualitatively matches the Eulerian statistics of τ_b . Despite the fluctuations, performance using the adaptive memory nears performance with the optimal memory (Figure 3C). This result confirms our intuition that memory should match the blank time. The advantage of the adaptive memory is that it relies solely on experience, with no prior information whatsoever. This is unlike T^* which can only be selected using prior information, with no guarantee of generalization to other plumes.

Learning to recover. So far, our agents combine a learned policy from non-void states to a heuristics from the void state, which we called the recovery strategy. We have considered a biologically-inspired heuristics where searchers make it back to locations within the plume by retracing their path backward. Similar results are obtained for Brownian recovery with a different optimal memory (see Materials and Methods and Supplementary Figure 1). To further strip the algorithm of heuristics, we ask whether the recovery strategy may be learned, rather than fixed a priori. To this end, we split the void state in many states, labeled with the time elapsed since first entering the void. The counter is reset to 0 whenever the searcher detects the odor. The definition of the 15 non-void states o_1, \dots, o_{15} remains unaltered. Interestingly, with this added degree of freedom, the agent learns an even better recovery strategy as reflected by all our measures of performance Figure 3D. Note that the learned recovery strategy resembles the casting behavior observed in flying insects [37], as discussed below.

Characterization of the optimal policies. To understand how different recoveries affect the agent’s behavior, we characterize the optimal policies obtained using the three recovery strategies. We visualize the probability to encounter each of the 16 olfactory states, or occupancy (circles in Figure 4), and the spatial distribution of the olfactory states. In the void state, the agent activates the recovery strategy. Recovery from the void state affects

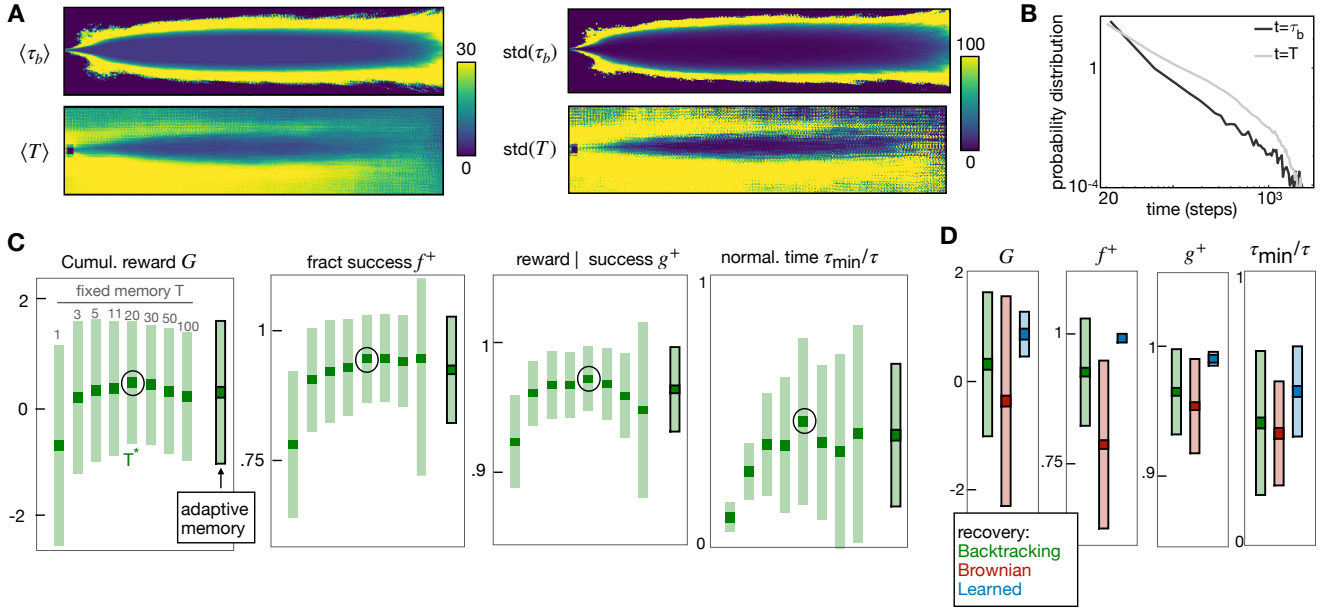


Figure 3: The adaptive memory approximates the duration of the blank dictated by physics and it is an efficient heuristics, especially when coupled with a learned recovery strategy. (A) Colormaps of the Eulerian blank times τ_b (top) and the sensing memory T (bottom): Left: averages; Right: standard deviations. The sensing memory statistics is computed over all agents that are located at each discrete cell, at any point in time. (B) Probability distribution of τ_b across all spatial locations and times (black) and of T across all agents at all times (gray). (C) Performance with the adaptive memory nears performance of the optimal fixed memory, here shown for backtracking; similar results apply to the Brownian recovery (Supplementary Figure 2). (D) Comparison of three recovery strategies with adaptive memory: The learned recovery with adaptive memory outperforms all fixed and adaptive memory agents. In (C) and (D) dark squares mark the mean, and light rectangles mark \pm standard deviation. No standard deviation is shown for the f^+ measure in the learned case as this strategy is deterministic (see Materials and Methods).

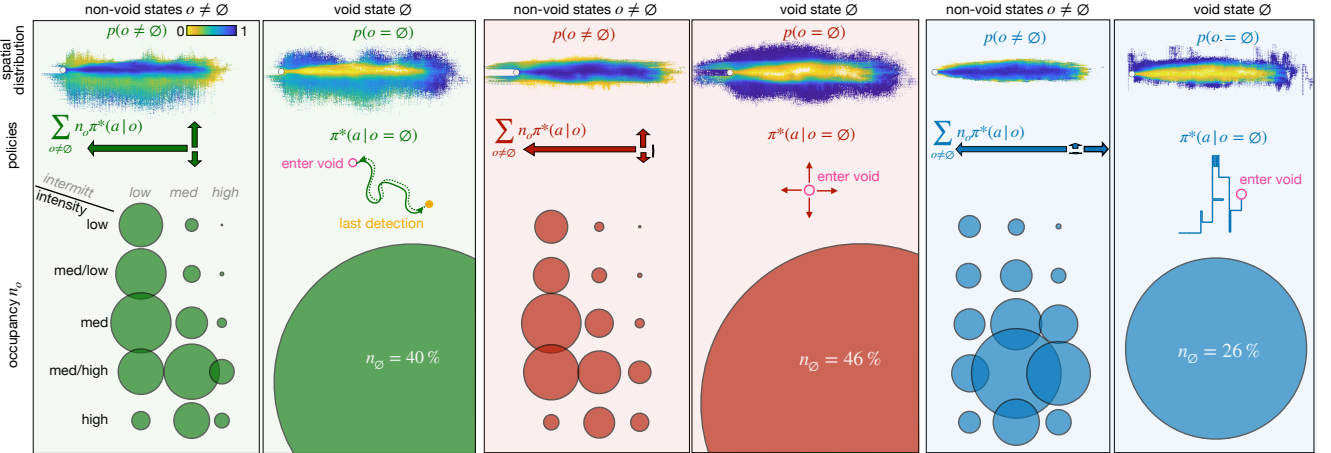


Figure 4: Optimal policies with adaptive memory for different recovery strategies: backtracking (green), Brownian (red) and learned (blue). For each recovery, we show the spatial distribution of the olfactory states (top); the policy (center) and the state occupancy (bottom) for non-void states (left) *vs* the void state $\pi^*(a|\emptyset)$ (right). Spatial distribution: probability that an agent at a given position is in any non-void olfactory state (left) or in the void state (right), color-coded from yellow to blue. Policy: actions learned in the non-void states $\sum_{o \neq \emptyset} n_o \pi^*(a|o)$, weighted on their occupancy n_o (left, arrows proportional to the frequency of the corresponding action) and schematic view of recovery policy in the void state (right). State occupancy: fraction of agents that is in any of the 15 non-void states (left) or in the void state (right) at any point in space and time. Occupancy is proportional to the radius of the corresponding circle. The position of the circle identifies the olfactory state (rows and columns indicate the discrete intensity and intermittency respectively). All statistics is computed over 43000 trajectories, starting from any location within the plume.

non-void olfactory states as well: their occupancy, their spatial distribution, and the action they elicit (Figure 4 and Supplementary Figure 3). This is because the agent computes its olfactory state online, according to its prior history which is affected by encounters with the void state. However, for all recoveries, non-void states are mostly encountered within the plume and largely elicit upwind motion (Figure 4, top, center). Thus macroscopically, all agents learn to surge upwind when they detect any odor within their memory, and to recover when their memory is empty. This suggests a considerable level of redundancy which may be leveraged to reduce the number of olfactory states, thus the computational cost. The void state shows the most relevant differences: for both heuristic recoveries, 40% or more of the agents are in the void state and they are spatially spread out. In contrast, in the case of learned recovery, the optimal policy limits occurrence of the void state to 26% of the agents, confined to a narrow band near the edge of the plume. From these locations, the agents quickly recover the plume, explaining the boost in performance discussed above. Note that, exclusively for the learned recovery, the optimal policy is enriched in actions down-

wind to avoid overshooting the source. Indeed, from positions beyond the source, the learned strategy is unable to recover the plume as it mostly casts sideways, with little to no downwind action. Intuitively, the precise locations where agents move downwind may be crucial to efficiently avoid overshooting. Thus the policy may depend on specific details of the odor plume, consistent with poorer generalization of the learned recovery (discussed next).

Tuning for adaptation to different environments.

Finally, we test performance of the trained agents on six environments, characterized by distinct fluid flows and odor plumes (Figure 5 and Materials and Methods). Environment 1 is the native environment, where the agents were originally trained; Environment 2 is obtained by increasing the threshold of detection, which makes the signals considerably more sparse with longer blanks. Environments 3 and 4 are closer to the lower surface of the simulated domain, where the plume is smaller and fluctuates less. Environment 5 is a similar geometry, but obtained for a smaller Reynolds number and a different way to generate turbulence. Finally Environment 6 has

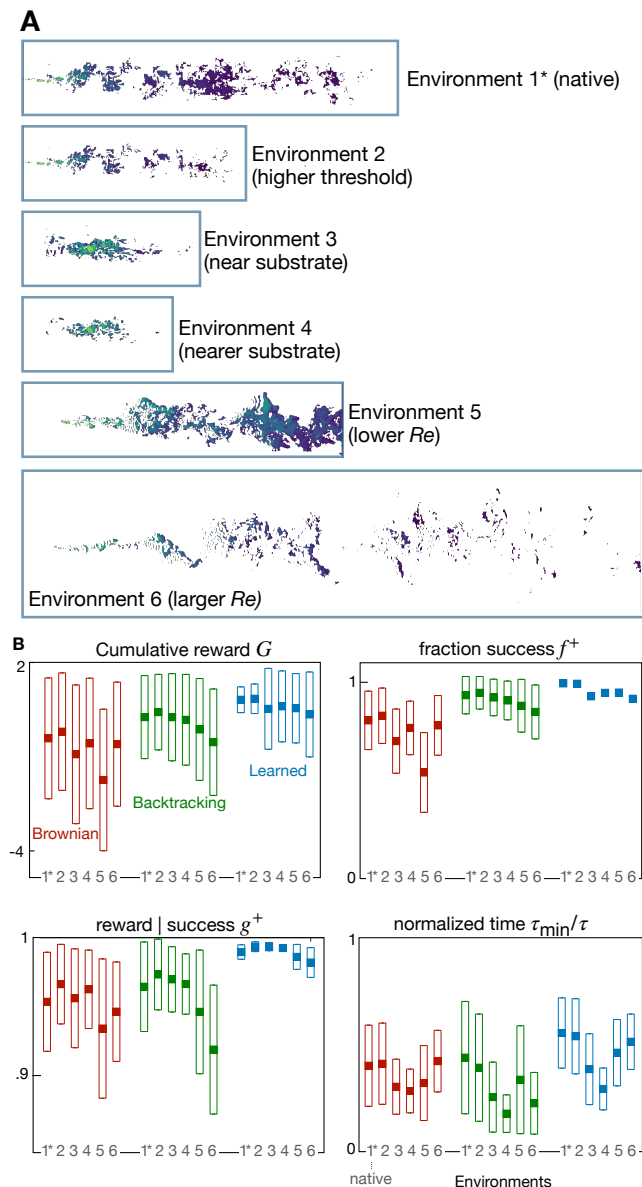


Figure 5: Generalization to statistically different environments. (A) Snapshots of odor concentration normalizes with concentration at the source, colorcoded from blue (0) to yellow (1) for environment 1 to 6 from top to bottom. Environment 1* is the native environment where all agents are trained. (B) Performance for the three recovery strategies backtracking (green) Brownian (red) and learned (blue), with adaptive memory, trained on the native environment and tested across all environments. Four measures of performance defined in the main text are shown. Dark squares mark the mean, and empty rectangles \pm standard deviation. No standard deviation is shown for the f^+ measure in the learned case as this strategy is deterministic (see Materials and Methods).

an even larger Reynolds number, a longer domain and a smaller source, which creates an even more dramatically sparse signal. We consider agents with adaptive memory and compare the three recovery strategies discussed above – Brownian, backtracking and learned, see Figure 5B. Comparing performance across environments we find that: (i) although performance is degraded when testing in non native environments, all agents with adaptive memory are still extremely likely to find the source; (ii) Brownian recovery has lowest performance and generalization across all environments; (iii) backtracking provides good performance and generalization; (iii) the learned recovery strategy performs best in all environments by all performance metrics. In the most intermittent Environment 6 a striking 91% of agents succeeds in finding the source, with trajectories less the twice as long as the shortest path to the source. The upshot of generalization is that agents may navigate distinct turbulent plumes using a baseline strategy learned in a specific plume. Importantly, even if performance (mildly) degrades, most agents still do reach the source, suggesting that fine-tuning this strategy may enable efficient adaptation to different environments. Further work is needed to establish how much fine-tuning is needed to fully adapt the baseline strategy to different environments.

Discussion

In this work, we showed that agents exposed to a turbulent plume learn to associate salient features of the odor time trace – the olfactory state – to an optimal move that guides them to the odor source. The upshot of responding solely to odor is that the agent does not navigate based on *where* it believes the target is and thus needs no map of space nor prior information about the odor plume, which avoids considerable computational burden. On the flip side, in our stimulus-response algorithm, agents need to start from within the plume, however sparse and fragmented. Indeed, far enough from the source, Q-learning agents are mostly in the void state and they can only recover the plume if they have previously detected the odor or are right outside the plume. In contrast, agents using a map of space can navigate from larger distances than are reachable by responding directly to odor cues. Indeed, in the map-based POMDP setting, absence of odor detection is still informative and it enables agents to first find the plume, and then refine the search to localize the target within the plume [26, 38].

We show that because turbulent odor plumes constantly

switch on and off, navigation must handle both absence and presence of odor stimuli. We address this fundamental issue by alternating between two distinct strategies: (i) Prolonged absence of odor prompts entry in the void state and triggers a recovery strategy to make contact with the plume again. We explored two heuristic recoveries and found that back-tracking to where the agent last detected odor is much more efficient than Brownian recovery. An even more efficient recovery can be learned that resembles cross-wind casting and limits the void state to a narrow region right outside of the plume. Casting is a well-studied computational strategy [12, 5] also observed in animal behavior, most famously in flying insects [37]. Intriguingly, cast and surge also emerges in algorithms making use of a model of the odor, whether for Bayesian updates or for policy optimization [3, 26, 30]. Whether natural casting behavior is learned, as in Q-learning, or is hard-wired in a model of the odor plume remains a fascinating question for further research. (ii) Odor detections prompt entry in non-void olfactory states, which predominantly elicit upwind surge. Blanks shorter than the sensing memory are ignored, i.e. agents do not enact recovery but respond to stimuli experienced prior to the short blank. Further work may optimize these non-void olfactory states by feature engineering, e.g. testing different discretizations to reduce redundancy or screening a large library of features using supervised learning as in [33] to potentially improve performance. Alternatively, feature engineering may be bypassed altogether by the use of recurrent neural networks (RNNs) as recently proposed in [29], possibly at the expense of interpretability. A systematic comparison using a common dataset is needed to elucidate how other heuristic and normative model-free algorithms handle odor presence *vs* odor absence.

To switch between the odor-driven strategy and the recovery strategy, we introduce a timescale T , which is an explicit form of temporal memory. T delimits a sensing window extending in the recent past, prior to the present time. All odor stimuli experienced within the sensing window affect the current response. By using fixed memories of different duration, we demonstrate that an optimal memory exists and that the optimal memory minimizes the occurrence of the void state. On the one hand, long memories are detrimental because they delay recovery from accidentally exiting the plume. On the other hand, short memories are detrimental because they trigger recovery unnecessarily, i.e. even for blanks typically experienced within the turbulent plume. The optimal memory thus matches the typical duration of the blanks. To avoid using prior

information on the statistics of the odor, we propose a simple heuristics setting memory adaptively equal to the most recent blank experienced along the path. The adaptive memory nears optimal performance despite dramatic fluctuations dictated by turbulence. Success of the heuristics suggests that a more accurate estimate of the future blank time may enable an even better adaptive memory; further work is needed to corroborate this idea.

Thus in Q-learning, memory is a temporal window matching odor blanks and distinguishing whether agents are in or out of the plume. The role of memory for olfactory search has been recently discussed in ref. [30]. In POMDPs, memory is stored in a detailed belief of agent position relative to the source. In finite state controllers, memory denotes an internal state of the agent and was linked to a coarse grained belief of the searcher being within or outside of the plume, similar to our findings. In recurrent neural networks memory is stored in the learned weights. A quantitative relationship between these different forms of memory and their connection to spatial perception remains to be understood.

We conclude by listing a series of experiments to test these ideas in living systems. First, olfactory search in living systems displays memory ([30, 10] and refs. therein). In insects, temporal scales can be measured associated to memory. Indeed, for flying insects loss of contact with a pheromone plume triggers crosswind casting and sometimes even downwind displacement [7, 39]. Interestingly, the onset of casting is delayed with respect to loss of contact with the plume, but this delay is not understood [39, 40]. In walking flies, the timing of previous odor encounters biases navigation [41]. (How) do these temporal timescales depend on the waiting times between previous detections? Using optogenetics [42, 43, 44, 45] or olfactory virtual reality with controlled odor delivery [46] experiments may measure memory as a function of the full history of odor traces. For insects, one may monitor memory by tracking the onset of cross wind casting with respect to the loss of the plume. More in general, a temporal memory may be defined by monitoring how far back in the past should two odor traces be identical in order to elicit the same repertoire of motor controls.

Second, our algorithm learns a stimulus-response strategy that relies solely on odor cues. The price to pay is that the agent must follow the ups and downs of the odor trace in order to compute averages and recognize blanks. A systematic study may use our algorithm to test the requirements of fidelity of this temporal representation, and how it depends on turbu-

lence. How does turbulence affect the fidelity of odor temporal representation in living systems? Crustaceans provide excellent model system to ask this question, as they are known to use bursting olfactory receptor neurons to encode temporal information from olfactory scenes [47, 48]. Temporal information is also encoded in the olfactory bulb of mammals [49, 50]. Organisms with chemo-tactile systems like the octopus [51] may serve as a comparative model, to ask whether touch-chemosensation displays a sloppier temporal response, reflecting that surface-bound stimuli are not intermittent.

Third, our Q-learning algorithm requires the agent to receive olfactory information, thus start near or within the odor plume. In contrast, algorithms making use of a spatial map and prior information on the odor plume may first search for the plume (in conditions of near zero information) and then search the target within the plume [26, 38, 3]. Animals are known to use prior information to home into regions of space where the target is more likely to be found; but they can switch to navigation in response to odor (see e.g. [7, 8, 9, 10]). What triggers the switch from navigation driven by spatial perception to navigation driven by odor? For mice, the need for spatial perception may be tested indirectly by comparing paths in light *vs* dark, noting that neuronal place fields, that mediate spatial perception, are better stabilized by vision than olfaction [52, 53]. Thus in the light, animals have the ability to implement both map-less and map-based algorithms, whereas in the dark they are expected to more heavily rely on map-less algorithms. To make sure animals start searching for the odor target even before sensing odor, operant conditioning can be deployed so that animals associate an external cue (e.g. a sound) to the beginning of the task.

The reinforcement learning view of olfactory navigation offers an exciting opportunity to probe how living systems interact with the environment to accomplish complex real-world tasks affected by uncertainty. Coupling time varying odor stimuli with spatial perception is an instance of the broader question asking how animals combine prior knowledge regarding the environment with reaction to sensory stimuli. We hope that our work will spark further progress into connecting these broader questions to the physics of fluids.

Acknowledgements

This research was supported by grants to A.S. from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation

Table 1: Data information

Feature	Value
Shape	1225×280
Number of time slices	2598
n_{thr}	3×10^{-6}
Source location	(20, 142)

programme (grant agreement No 101002724 RIDING), the Air Force Office of Scientific Research under award number FA8655-20-1-7028, and the National Institutes of Health (NIH) under award number R01DC018789. L.R. acknowledges the financial support of the European Research Council (grant SLING 819789). We thank Francesco Viola for support and discussions regarding computational fluid dynamics as well as Antonio Celani, Venkatesh Murthy, Yujia Qi, Francesco Boccardo, Luca Gagliardi, Francesco Marcolli and Arnaud Ruymaekers for comments on the manuscript.

Materials and Methods

Data description. The data we used to train the agents is a set of 2598 matrices $\{D_t\}_{t=1}^{2598}$. Every matrix $D_t \in \mathbb{R}^{1225 \times 280}$ contains the odor intensity in every position (i, j) i.e. $(D_t)_{i,j}$ represents the odor intensity in position (i, j) at time t . The source of odor is in position (20, 142) and, in order to simplify the training, we considered as terminal states every position in a circle centered in the source position and with radius 10 called *source region*. Data are obtained from a direct numerical simulation of the Navier-Stokes equations and the equations of transport of the odor. Environments 1 to 4 are derived from simulations of a channel flow described in ref [33]; Environment 5 corresponds to an additional simulation described below. We preprocess the data to eliminate simulation noise by setting to zero every entry of these matrices when they are smaller than a *noise level* $n_{\text{lvl}} := 3 \times 10^{-6}$. Data information are summarized in Table 1.

In environment 5, the odor is advected by a turbulent open channel flow, with three hemispherical obstacles placed on the ground close to the inlet to generate turbulence. The Navier-Stokes equations (1) and advection-diffusion equation for odor transport (2) are solved using a central second order finite difference scheme. The convective terms are discretized in time using an explicit Adams – Bashforth method; and the viscous and diffusion terms using an implicit Crank-Nicolson method [54]. The code is written in Fortran and is GPU parallelized. The channel is divided into $1024 \times 256 \times 128$ grid points along streamwise, spanwise and wall-normal directions

respectively. The corresponding average spatial resolutions are $\Delta x = 5\eta, \Delta y = 5\eta, \Delta z = 4\eta$, where η is the Kolmogorov length scale. Three hemispheres of radius 100η are placed at a distance of 250η from the inlet on the ground, equally spaced along the spanwise direction. The channel is forced using a constant pressure gradient. For the velocity field, we impose a no-slip boundary condition at the ground and on the obstacles ($\mathbf{u} = 0$) and a free-slip boundary on top ($u_z = 0, \partial_z u_x = \partial_z u_y = 0$). The velocity field is periodic along the streamwise and spanwise directions. The bulk Reynolds number is 7800. For the odor field, we impose Dirichlet condition ($c = 0$) at the ground, on the obstacles and inlet, no-flux ($\partial_z c = 0$) on top, and outflow along other directions. Similar to the native environment, we choose the Schmidt number to be 1. The odor source is located downstream of the obstacle and centered at $[640\eta, 640\eta, 256\eta]$ along streamwise, spanwise and wall-normal directions respectively. The odor source has a Gaussian profile with a standard deviation of 8η .

Environment 6 is similar to environment 5 albeit with a higher bulk Reynolds number of 17500. Here, the channel is divided into $2000 \times 500 \times 200$ grid points and has an average spatial resolution of $\Delta x = \Delta y = \Delta z = 5.5\eta$. The odor source has a Gaussian profile centered at $[825\eta, 1375\eta, 550\eta]$ with a standard deviation of 3η .

$$\rho \left(\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} \right) = -\nabla P + \mu \nabla^2 \mathbf{u} + \mathbf{f}; \quad (1)$$

$$\nabla \cdot \mathbf{u} = 0.$$

$$\frac{\partial c}{\partial t} + \mathbf{u} \cdot \nabla c = D \nabla^2 c + s. \quad (2)$$

Olfactory states, Features & Discretization.

Each agent stores the odor concentrations detected in the previous T time steps in a vector \mathbf{M} . We introduce an adaptive sensitivity threshold function $s_{\text{thr}}(\cdot)$ defined as

$$s_{\text{thr}}(T) := \max \left\{ \frac{C_{\text{thr}}}{T} \sum_{i=1}^T M_i, n_{\text{thr}} \right\}, \quad (3)$$

where M_i denotes the i -th element of M and $C_{\text{thr}} > 0$ is a scaling constant (in our experiments we set it as 0.5). T denotes the cardinality of M . Given a memory M , we can define the filtered memory Δ^M as the set which contains every element of the memory M that is higher than the sensitivity threshold $s_{\text{thr}}(M)$ i.e.

$$\Delta^M := \{z \in M \mid z > s_{\text{thr}}(M)\}. \quad (4)$$

Then at timestep t , given the agent memory M_t , we define the average intensity $c(M_t)$ and the intermittency

$i(M_t)$ as:

$$c(M_t) := \begin{cases} \frac{1}{|\Delta^{M_t}|} \sum_{i=1}^{|\Delta^{M_t}|} (\Delta^{M_t})_i, & |\Delta^{M_t}| > 0 \\ 0 & \end{cases}, \quad (5)$$

$$i(M_t) := \frac{|\Delta^{M_t}|}{|M_t|}.$$

Note that the average intensity is defined on the filtered memory Δ^M , i.e. conditioned to detecting odors above threshold. Since the features defined in (5) returns real numbers, in order to use (tabular) q-learning, we need to discretize them. We denote with $\bar{i}(M_t)$ the discretized intermittency. This is defined as follow

$$\bar{i}(M_t) := \begin{cases} 0, & \text{if } i(M_t) \leq 0.33 \\ 1, & \text{if } 0.33 < i(M_t) \leq 0.66 \\ 2, & \text{if } i(M_t) > 0.66 \end{cases}. \quad (6)$$

The average intensity is bounded between zero and the maximum concentration of odor at the source. To avoid prior information on the source, we use a more structured procedure to discretize the average intensity online, based on the agent's experience only. At every timestep t , the average intensity $c(M_t)$ is computed and collected in a dataset X_t i.e.

$$X_t := \{c(M_0), \dots, c(M_t)\}.$$

Then, its discretized value is obtained by the following rule:

$$\bar{c}(M_t, X_t) := \begin{cases} 0, & c(M_t) \leq p(X_t, 25) \\ 1, & p(X_t, 25) < c(M_t) \leq p(X_t, 50) \\ 2, & p(X_t, 50) < c(M_t) \leq p(X_t, 80) \\ 3, & p(X_t, 80) < c(M_t) \leq p(X_t, 99) \\ 4, & c(M_t) > p(X_t, 99) \end{cases}, \quad (7)$$

where $p(X_t, n)$ denotes the n -th percentile of X_t . Finally, we can define the feature map ϕ_t as a function of the memory M_t and the dataset of average intensities X_t at timestep t

$$\phi_t(M_t, X_t) := [\bar{i}(M_t), \bar{c}(M_t, X_t)].$$

This defines the current olfactory state s_t i.e. at timestep t , the agent is in the olfactory state $o_t := \phi_t(M_t, X_t)$. The case where the agent has no odor detections above threshold in its current memory, i.e. $|\Delta(M_t)| = 0$ corresponds to an additional state called void state (\emptyset) in the main text.

Agent Behavior and Policies. Now, we describe how the agent interacts with the environment to solve the navigation problem. At every time step $t \in \mathbb{N}$, the agent observes an odor point z_t and updates its memory including the new observation and removing the oldest i.e. it defines a memory M_t with the following rule

$$M_t := \left[\left(M_{t-1} \right)_2, \dots, \left(M_{t-1} \right)_{|M_{t-1}|}, o_t \right]. \quad (8)$$

Then, it updates the dataset of average intensities i.e. $X_t := X_{t-1} \cup \{c(M_t)\}$ and it computes the olfactory state o_t . According to o_t , the agent chooses an action a_t using a policy. As indicated in the main text, actions are the coordinate directions i.e. we define an action set \mathcal{A} as follow

$$\mathcal{A} := \{e_1, e_2, -e_1, -e_2\},$$

where e_i denotes the i -th canonical base. As explained in the main text, actions are selected using one of two policies according to the current olfactory state o_t . More precisely, if the olfactory state o_t is not the void state, then the (ϵ -greedy) Q-learning policy is used. Formally, let Q be the Q matrix of the agent and let $o_t \neq \emptyset$, then the agent plays the action a_t such that

$$a_t = \begin{cases} a \in \arg \max_{a \in \mathcal{A}} Q(o_t, a) & \text{with probability } 1 - \epsilon \\ a \sim \mathcal{U}(\mathcal{A}) & \text{with probability } \epsilon \end{cases}, \quad (9)$$

where, with $a \sim \mathcal{U}(\mathcal{A})$, we indicate an action a uniformly sampled from \mathcal{A} . At test phase, the exploration-exploitation parameter ϵ is set to 0 and, thus, in an olfactory state $o_t \neq \emptyset$ the policy is deterministic. While training phase behavior is described in next paragraphs. In the void state $o_t = \emptyset$, the agent chooses the action $a_t \in \mathcal{A}$ according to a separated policy called *recovery strategy*. In our experiments, we defined and compared three different recovery strategies: Brownian, Backtracking and Learned.

Brownian recovery. It is the simplest strategy we consider, consisting of playing random actions in the void state. Suppose that at time step t , the agent is in the void olfactory state, i.e., $o_t = \emptyset$, then a_t is sampled uniformly from the action set \mathcal{A} . However, it is important to note for long memories agents start to recover when they are already far from the plume, and hitting the plume by random walk is prohibitively long. To avoid wandering away from the plume, the memory is constrained to be shorter, consistent with the observation that the optimal memory is $T^* = 3$ to 5, much shorter than for backtracking. At this memory, several blanks within the plume will cause the agent to recover, hence the lower performance of the Brownian recovery.

Backtracking Recovery. In order to accelerate recovery from accidentally exiting the plume, we let the agents backtrack to the position where they last detected the odor. To this end, we first enumerate the actions with numbers from one to four. Then we introduce a new memory called *action memory* A . For simplicity, we consider the setting in which $|A| = |M|$. At time-step $t = 0$, this memory is initialized as a vector of zeros indicating that the action memory is empty i.e. we define $A_0 \in \mathbb{N}^{|M|}$ such that for every $i = 1, \dots, |A|$

$$A_i = 0.$$

For every timestep $t > 0$, the agent observes an odor point z_t and updates the memory through (8). Moreover, the action memory is updated according to the status of the memory. If the last observation is smaller than the sensitivity threshold i.e. $z_t < s_{\text{thr}}(M_t)$, the action previously played a_{t-1} (represented by a natural number in $[1, 4]$) is stored in the action memory i.e. for some $\Delta > 0$, let

$$A_{t-1} = [a_{t-\Delta}, \dots, a_{t-2}, 0, \dots, 0].$$

Then

$$A_t = [a_{t-\Delta}, \dots, a_{t-2}, a_{t-1}, \dots, 0].$$

If at time-step t , the observation z_t is larger than the sensitivity threshold then the action memory is reset i.e. $A_t \in \mathbb{N}^{|M|}$ with $(A_t)_i = 0$ for every i . If at timestep t , the memory is empty i.e. $c(M_t) = 0$, then the backtracking procedure is executed: the last non-zero element of the action memory is extracted and the inverse action is played i.e. For some $\Delta > 0$, let

$$A_{t-1} = [a_{t-\Delta}, \dots, a_{t-2}].$$

Then, it plays the action a_{t-2} and updates the action memory as follow

$$A_t = [a_{t-\Delta}, \dots, a_{t-3}, 0].$$

This procedure is repeated until either an observation larger than the sensitivity threshold is obtained or the action memory becomes empty. In the former case, the action memory is cleared and the action is chosen according to the Q-learning policy ((9)). In the latter case, a random action is played.

Note that this strategy only provides exploration after the backtracking fails to recover detections. Also, if agents start with no detection at time 0, the procedure is equivalent to Brownian motion.

Learned recovery. In this case recovery policy is learned by splitting the void state in several states labeled by the time since entry in the void state. In our

experiments, we split the void state in 30 states. Actions are then learned as in all other non-void states and the optimal action is always chosen with (9).

Training An agent start at a random location within the odor plume at time 0. Its memory is initialized with the prior $|M_0|$ odor detections at its initial location $M_0 = [z_{-|M_0|}, \dots, z_0]$, obtained from the fluid dynamics simulation. The Q-function Q_0 is initialized with 0.6 for all actions and olfactory states. The first dataset of average intensities contains the first value $X_0 = \{c(M_0)\}$. At every times step $t > 0$, the agent gets an odor observation z_t from its new position and updates its memory including the new observation and removing the oldest and the olfactory state o_t is computed (as described in previous paragraphs). The dataset of average intensities is updated: $X_t = X_{t-1} \cup \{c(M_t)\}$. Exploration-exploitation parameter ϵ_k is scheduled as follow

$$\epsilon_k = \eta_{\text{init}} \exp(-\eta_{\text{decay}}k),$$

where, in our experiments, $\eta_{\text{init}} = 0.99$ and $\eta_{\text{decay}} = 0.0001$. At every episode k , the Q-function is updated at every time step t as

$$Q_{k+1}(s_t, a_t) := (1 - \alpha_k)Q_k(s_t, a_t) + \alpha_k(r_t + \gamma \max_{a'} Q_k(s_{t+1}, a')),$$

where R_t is the immediate reward received playing the action a_t . o_t and o_{t+1} are the current and the next olfactory states and α_k is the learning rate at episode k . This is scheduled as

$$\alpha_k = \alpha_{\text{init}} \exp(-\alpha_{\text{decay}}k),$$

where, in our experiments, $\alpha_{\text{init}} = 0.25$ and $\alpha_{\text{decay}} = 0.001$. For the experiments, agents are trained in 100000 episodes and an horizon of 5000 steps. The agent velocity is set to 10 and the discount factor is $\gamma = 0.9999$.

Agents Evaluation. To evaluate the performance of the different agents, we consider four metrics: the cumulative reward G (which is the actual quantity that the algorithm optimizes for); normalized time (defined below); the fraction of success f^+ and the value conditioned on success g^+ . For a fixed position (i, j) , we denote with $\tau_{\text{min}}(i, j)$ the minimum number of steps required to reach the source region from (i, j) i.e. the length of the shortest path.

We define D_{init} the set of points in which the first observation is above the sensitivity threshold (valid points). For each initial position $(i, j) \in D_{\text{init}}$, let $\tau(i, j)$ be the duration of the path obtained by an agent to reach the source. Note that $\tau(i, j)$ is a random variable for the stochastic backtracking and Brownian recoveries, but it

is deterministic for the learned strategy that has no random components. For each admissible location (i, j) , we define four performance metrics:

$$\begin{aligned} G(i, j) &= \langle e^{-\lambda\tau(i, j)} - \frac{\eta}{1-\gamma}(1 - e^{-\lambda\tau(i, j)}) \rangle \\ f^+(i, j) &= \frac{n_{\text{success}}(i, j)}{n_{\text{reps}}} \\ g^+(i, j) &= \langle e^{-\lambda\tau(i, j)} | \text{success} \rangle \\ \frac{\tau_{\text{min}}}{\tau}(i, j) &= \langle \frac{\tau_{\text{min}}(i, j)}{\tau(i, j)} \rangle \end{aligned}$$

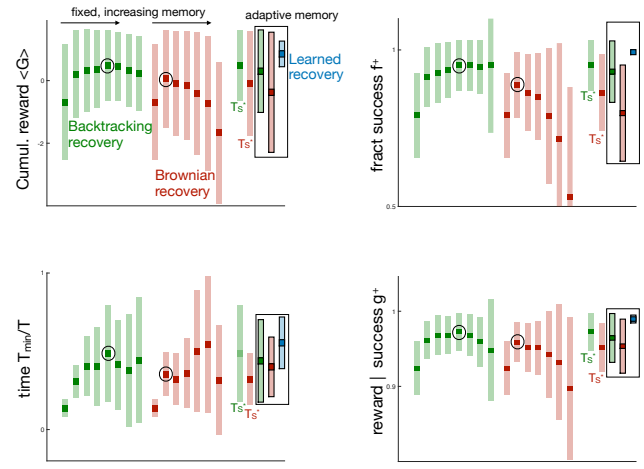
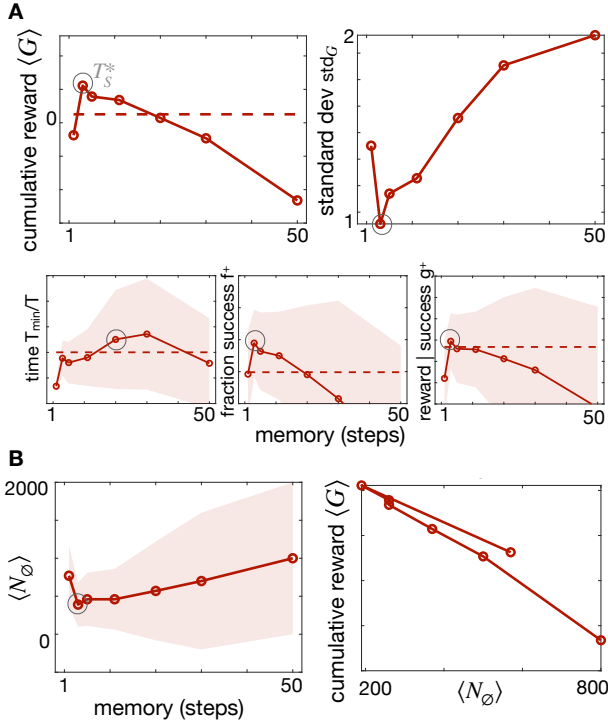
where n_{reps} is the number of test trajectories from each admissible location, and we use $n_{\text{reps}} = 10$. We then compute statistics of the performance metrics over the D_{init} initial positions and report the average ($\langle \cdot \rangle$) and standard deviation (std). Note that for the learned strategy, $\tau(i, j)$ is deterministic, hence $f^+(i, j)$ is 0 or 1 and therefore we omit its standard deviation.

References

- [1] H C Berg. Chemotaxis in bacteria. *Annual Review of Biophysics and Bioengineering*, 4(1):119–136, 1975.
- [2] J. Murlis, J. S. Elkinton, and R. T. Cardé. Odor plumes and how insects use them. *Annual Review of Entomology*, 37:505, 1992.
- [3] M. Vergassola, E. Villermaux, and B.I. Shraiman. 'Infotaxis' as a strategy for searching without gradients. *Nature*, 445:406, 2007.
- [4] B. Shraiman and E. Siggia. Scalar turbulence. *Nature*, 405:639, 2000.
- [5] E Balkovsky and Shraiman B. I. Olfactory search at high reynolds number. *Proc Nat Acad Sci*, 99(20):12589–93, 2002.
- [6] Gautam Reddy, Venkatesh N. Murthy, and Massimo Vergassola. Olfactory sensing and navigation in turbulent environments. *Annual Review of Condensed Matter Physics*, 13(1):191–213, 2022.
- [7] R. T. Cardé. Navigation along windborne plumes of pheromone and resource-linked odors. *Annual Review of Entomology*, 66:317, 2021.
- [8] C. Schal. Intraspecific vertical stratification as a mate-finding mechanism in tropical cockroaches. *Science*, 215:1505, 1982.
- [9] DH. Gire, V. Kapoor, A. Arrighi-Allisan, A. Seminara, and VN. Murthy. Mice develop efficient strategies for foraging and navigation using complex natural stimuli. *Curr Biol*, 26:1261, 2016.

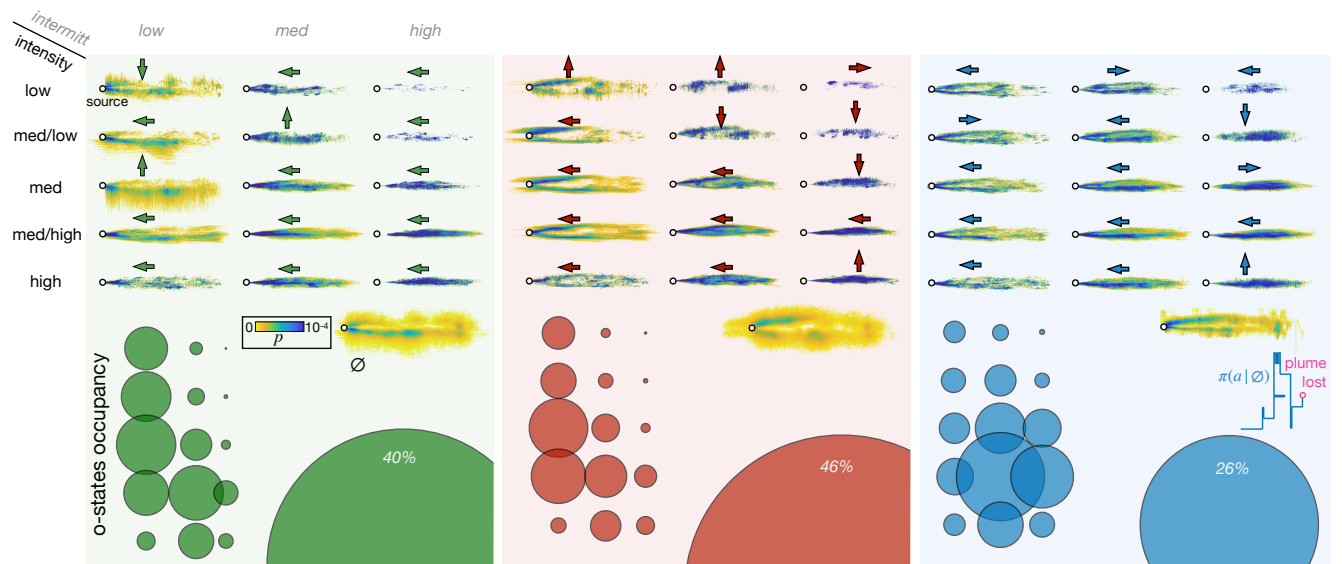
- [10] K. L. Baker, M. Dickinson, T. M. Findley, D. H. Gire, M. Louis, M. P. Suver, J. V. Verhagen, K. I. Nagel, and M. C. Smear. Algorithms for olfactory search across species. *Journal of Neuroscience*, 38:9383, 2018.
- [11] A. Celani and E. Panizon. Olfactory search. *in review*, .:., 2024.
- [12] Baker T. C. Upwind flight and casting flight: complementary and tonic systems used for location of sex pheromone sources by male moths. *Proc. 10th Intl Symposium on Olfaction and Taste*, 13:18, 1990.
- [13] E. Kramer. A tentative intercausal nexus and its computer model on insect orientation in windborne pheromone plumes. *in Insect Pher. Res, New Dir.*, .:232, 1997.
- [14] J.H. Belanger and M.A. Willis. Biologically-inspired search algorithms for locating unseen odor sources. *In, Proc. IEEE Symp. Intell. Control (ISIC '98) and IEEE Symp. Comp. Intell. Robot. Autom. (CIRA '98)*, .:265, 1988.
- [15] J Atema. Eddy chemotaxis and odor landscapes: exploration of nature with animal sensors. *Biol. Bull.*, 191:129, 1996.
- [16] BT Michaelis, KW Leathers, YV Bobkov, BW Ache, JC Principe, R Baharloo, IM Park, and MA Reidenbach. Odor tracking in aquatic organisms: the importance of temporal and spatial intermittency of the turbulent plume. *Sci. Rep.*, 10:7961, 2020.
- [17] Mahmut Demir, Nirag Kadakia, Hope D Anderson, Damon A Clark, and Thierry Emonet. Walking *Drosophila* navigate complex plumes using stochastic decisions biased by the timing of odor encounters. *eLife*, 9:e57524, 2020.
- [18] R. Sutton and A. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- [19] A. R. Cassandra, L. P. Kaelbling, and J. A. Kurien. Acting under uncertainty: Discrete bayesian models for mobile-robot navigation. *Proc IEEE/RSJ Internl Conf Intelligent Robots and Systems. IROS '96*, 2:963, 1996.
- [20] Steven M LaValle. *Planning algorithms*. Cambridge University Press, 2006.
- [21] A. Loisy and C. Eloy. Searching for a source without gradients, how good is infotaxis and how to beat it. *Proc. R. Soc. A*, 478:20220118, 2022.
- [22] H. Ishida. Chemical sensing in robotic applications: a review. *IEEE Sensors, J.*, 12:3163, 2020.
- [23] V Krishnamurthy. *Partially Observed Markov Decision Processes*. Cambridge University Press, 2016.
- [24] M Hauskrecht. Value-function approximations for partially observable markov decision processes. *J. Artif. Intell. Res.*, 13:33, 2000.
- [25] G Shani, J Pineau, and R Kaplow. A survey of point-based pomdp solvers. *Autonomous Agents and MultiAgent Systems*, 27:1–51, 2013.
- [26] Nicola Rigolli, Gautam Reddy, Agnese Seminara, and Massimo Vergassola. Alternation emerges as a multi-modal strategy for turbulent odor navigation. *eLife*, 11:e76989, aug 2022.
- [27] R. A. Heinonen, L. Biferale, A. Celani, and M. Vergassola. Optimal policies for bayesian olfactory search in turbulent flows. *Phys. Rev. E*, 107:055105, May 2023.
- [28] Aurore Loisy and Robin A. Heinonen. Deep reinforcement learning for the olfactory search pomdp: a quantitative benchmark. *Cereb CortexThe European Physical Journal E*, 46:17, 2023.
- [29] Satpreet H. Singh, Floris van Breugel, Rajesh P. N. Rao, and Bingni W. Brunton. Emergent behaviour and neural dynamics in artificial agents tracking odour plumes. *Nature Machine Intelligence*, 5:58–70, 2023.
- [30] K. V. B. Verano, E Panizon, and A Celani. Olfactory search with finite-state controllers. *Proc Nat Acad Sci*, 120(34):e2304230120, 2023.
- [31] G. Falkovich, K. Gawedzki, and M. Vergassola. Particles and fields in fluid turbulence. *Rev. Mod. Phys.*, 73:913, 2001.
- [32] A. Celani, E. Villermaux, and M. Vergassola. Odor landscapes in turbulent environments. *Phys. Rev. X*, 4:041015, 2014.
- [33] Nicola Rigolli, Nicodemo Magnoli, Lorenzo Rosasco, and Agnese Seminara. Learning to predict target location with turbulent odor plumes. *eLife*, 11:e72196, aug 2022.
- [34] Ariane S. Etienne and Kathryn J. Jeffery. Path integration in mammals. *Hippocampus*, 14(2):180–192, 2004.

- [35] Ariane S. Etienne, Roland Maurer, and Valérie Séguinot. Path Integration in Mammals and its Interaction With Visual Landmarks. *Journal of Experimental Biology*, 199(1):201–209, 01 1996.
- [36] S. Heinze, A. Narendra, and A. Cheung. Principles of insect path integration. *Current Biology*, 28:R1043, 2018.
- [37] C. T. David, J. S. Kennedy, and A. R. Ludlow. Finding of a sex pheromone source by gypsy moths released in the field. *Nature*, 303:804–806, 1983.
- [38] A. Loisy and R. A. Heinonen. Deep reinforcement learning for the olfactory search pomdp: a quantitative benchmark. *European Physical Journal E*, 46:17, 2023.
- [39] L. P. S. Kuenen and R. T. Cardé. Strategies for recontacting a los pheromone plume: casting and upwind flight in the male gypsy moth. *Physiological Entomology*, 15:317, 1994.
- [40] van Breugel F and Dickinson MH. Plume-tracking behavior of flying drosophila emerges from a set of distinct sensory-motor reflexes. *Curr Biol*, 24:274, 2014.
- [41] Mahmut Demir, Nirag Kadakia, Hope D Anderson, Damon A Clark, and Thierry Emonet. Walking *Drosophila* navigate complex plumes using stochastic decisions biased by the timing of odor encounters. *eLife*, 9:e57524, nov 2020.
- [42] Ruben Gepner, Mirna Mihovilovic Skanata, Natalie M Bernat, Margarita Kaplow, and Marc Gershow. Computations underlying *Drosophila* phototaxis, odor-taxis, and multi-sensory integration. *eLife*, 4:e06229, may 2015.
- [43] Luis Hernandez-Nunez, Jonas Belina, Mason Klein, Guangwei Si, Lindsey Claus, John R Carlson, and Aravinthan DT Samuel. Reverse-correlation analysis of navigation dynamics in *Drosophila* larva using optogenetics. *eLife*, 4:e06225, may 2015.
- [44] Andrew M. M. Matheson, Aaron J. Lanz, Ashley M. Medina, Al M. Licata, Timothy A. Currier, Mubarak H. Syed, and Katherine I. Nagel. A neural circuit for wind-guided olfactory navigation. *Nature Communications*, 13:4613, 2022.
- [45] S. D. David Stupski and F. van Breugel. Wind gates search states in free flight. *bioArX*, doi.org/10.1101/2023.11.30.569086:1, 2024.
- [46] Brad A. Radvansky and Daniel A. Dombeck. An olfactory virtual reality system for mice. *Nature Communications*, 9:839, 2018.
- [47] Y.V. Bobkov and B. Ache. Intrinsically bursting olfactory receptor neurons. *J. Neurophysiol*, 97:1052, 2007.
- [48] B.W. Ache, A.M. Hein, Y.V. Bobkov, and J.C. Principe. Smelling time: A neural basis for olfactory scene analysis. *Trends Neurosci.*, 39:649–655, 2016.
- [49] Ryan M. Carey, Justus V. Verhagen, Daniel W. Wesson, Nicolás Pérez, and Matt Wachowiak. Temporal structure of receptor neuron input to the olfactory bulb imaged in behaving rats. *Journal of Neurophysiology*, 101(2):1073–1088, 2009.
- [50] T. Ackels, A. Erskine, and D. et al. Dasgupta. Fast odour dynamics are encoded in the olfactory system and guide behaviour. *Nature*, 593:558, 2021.
- [51] CAH Allard, G Kang, JJ Kim, WA Valencia-Montoya, RE Hibbs, and NW Bellono. Structural basis of sensory receptor evolution in octopus. *Nature*, 616:373, 2023.
- [52] E. Save, L. Nerad, and B. Poucet. Contribution of multiple sensory information to place field stability in hippocampal place cells. *Hippocampus*, 10:64, 2000.
- [53] S Zhang and Manahan-Vaughan D. Spatial olfactory learning contributes to place field formation in the hippocampus. *Cereb Cortex*, 25:423–32, 2015.
- [54] Francesco Viola, Valentina Meschini, and Roberto Verzicco. Fluid–structure–electrophysiology interaction (fsei) in the left-heart: a multi-way coupled computational model. *European Journal of Mechanics-B/Fluids*, 79:212–232, 2020.



Supplementary Figure 1 : The role of temporal memory with Brownian recovery strategy (same as main Figure 2A). Total cumulative reward (top left) and standard deviation (top right) as a function of memory showing an optimal memory $T^* = 3$ for the Brownian agent. Other measures of performance with their standard deviations show the same optimal memory (bottom). The tradeoff between long and short memories discussed in the main text holds, but here exiting the plume is much more detrimental because regaining position within the plume by Brownian motion is much lengthier.

Supplementary Figure 2 : All four measures of performance across all agents with fixed and adaptive memory and with adaptive memory for the three recovery strategies.



Supplementary Figure 3 : Optimal policies for different recovery strategies and adaptive memory. From left to right: results for backtracking (green), Brownian (red) and learned (blue) recovery strategies. Top: probability that an agent in a given olfactory state is at a specific spatial location color-coded from yellow to blue. Rows and columns indicate the olfactory state; the void state is in the lower right corner. Arrows indicate the optimal action from that state. Bottom: Circles represent occupancy of each state, olfactory states are arranged as in the top panel. All statistics is computed over 43000 trajectories, starting from any location within the plume.