# Neural optimal feedback control
# with local learning rules

**Johannes Friedrich** [1]    **Siavash Golkar** [1]    **Shiva Farashahi** [1]

**Alexander Genkin** [5]    **Anirvan M. Sengupta** [2,3,4]    **Dmitri B. Chklovskii** [1,5]

[1] Center for Computational Neuroscience, Flatiron Institute
[2] Center for Computational Mathematics, Flatiron Institute
[3] Center for Computational Quantum Physics, Flatiron Institute
[4] Department of Physics and Astronomy, Rutgers University
[5] Neuroscience Institute, NYU Medical Center

`{jfriedrich,sgolkar,sfarashahi,dchklovskii}@flatironinstitute.org`
`{alexander.genkin,anirvans.physics}@gmail.com`

## Abstract

A major problem in motor control is understanding how the brain plans and executes proper movements in the face of delayed and noisy stimuli. A prominent framework for addressing such control problems is Optimal Feedback Control (OFC). OFC generates control actions that optimize behaviorally relevant criteria by integrating noisy sensory stimuli and the predictions of an internal model using the Kalman filter or its extensions. However, a satisfactory neural model of Kalman filtering and control is lacking because existing proposals have the following limitations: not considering the delay of sensory feedback, training in alternating phases, and requiring knowledge of the noise covariance matrices, as well as that of systems dynamics. Moreover, the majority of these studies considered Kalman filtering in isolation, and not jointly with control. To address these shortcomings, we introduce a novel online algorithm which combines adaptive Kalman filtering with a model free control approach (i.e., policy gradient algorithm). We implement this algorithm in a biologically plausible neural network with local synaptic plasticity rules. This network performs system identification and Kalman filtering, without the need for multiple phases with distinct update rules or the knowledge of the noise covariances. It can perform state estimation with delayed sensory feedback, with the help of an internal model. It learns the control policy without requiring any knowledge of the dynamics, thus avoiding the need for weight transport. In this way, our implementation of OFC solves the credit assignment problem needed to produce the appropriate sensory-motor control in the presence of stimulus delay.

## 1 Introduction

The sensorimotor control system has exceptional abilities to perform fast and accurate movements in a variety of situations. To achieve such skillful control, this system faces two key challenges: (i) sensory stimuli are noisy, making estimation of current state of the system difficult, and (ii) sensory stimuli is often delayed, which if unaccounted, results in movements that are inaccurate and unstable [1]. Optimal Feedback Control (OFC) has been proposed as a solution to this control problem [2, 3]. OFC approaches these problems by building an internal model of the system dynamics, and using

this internal model to generate control actions. OFC often employs Kalman filtering to optimally integrate the predictions of this internal model and the noisy/delayed sensory stimuli.

Because of the power and flexibility of the OFC framework, biologically plausible neural architectures capable of building such internal models has been under active investigation. Specifically, earlier works used attractor dynamics implemented through a recurrent basis function network [4] or a line attractor network [5] to implement Kalman filters. Kalman filtering and control has also been implemented through different phases of estimation, system identification and control [6], and more recently, using a particle filtering method for Kalman filtering [7].

Nonetheless, these works suffer from major limitations. Importantly, none of these considered that sensory feedback is delayed [4, 6, 5, 7, 8], although it has been prominent in the original computational-level OFC proposal [2], or merely considered the case of Kalman filtering, and not the combination of it with control [4, 5, 7, 8]. These works also required knowledge of the noise covariances, either a priori [4, 5, 7, 8] or obtained in a separate 'offline sensor' mode [6]. Moreover, many of these works lack biological plausibility and realism one would expect from a viable model of brain function [4, 5, 7, 8]. Crucially, biological plausibility requires the network to operate online, (i.e. receive a stream of noisy measurement data and process them on the fly), and also requires synaptic plasticity rules to be local (i.e. learn using rules that only depend on variables represented in pre and postsynaptic neurons and/or on global neuromodulatory signals). Lastly, several of these models suffer from combinatorial explosion as the dimensionality of the input grows [4, 5], require running an inner loop until convergence at each time step [8, 6], or require separate learning and execution phases [6], cf. Table 1.

We address these shortcomings and present a complete neural implementation of optimal feedback control, thus tackling an open issue in biological control [9]. In this model, which we call Bio-OFC, the state space, the prediction error [10, 11] (i.e., the mismatch between the network's internal prediction and delayed sensory feedback), and the control are represented by different neurons, cf. Fig. 1. The network also receives scalar feedback related to the objective function, as a global signal, and utilizes this signal to update the synaptic connection according to policy gradient method [12, 13]. To test the performance of our network, we simulate Bio-OFC in episodic (finite horizon) tasks (e.g., a discrete-time double integrator model, a hand reaching task [1], and a simplified fly simulation).

**Summary of contributions:**

- We introduce Bio-OFC, a biologically plausible neural network that combines adaptive model based state discovery via adaptive Kalman filtering with a model free control agent.

- Our implementation does not require knowledge of noise covariances nor the system dynamic, considers delayed sensory feedback, and has no separate learning/execution phases.

- Our model-free control agent enables closed-loop control, thus avoiding the weight transport problem, a challenging problem even in non-biological control. [14, 15]

Table 1: **Limitations of previously proposed neural implementations of OFC.** Presence or absence of different properties in previously proposed neural models, and their comparison to Bio-OFC. Guide to symbols: ✓: true, ✗: false, ✓✗: partially true, N/A: not applicable.

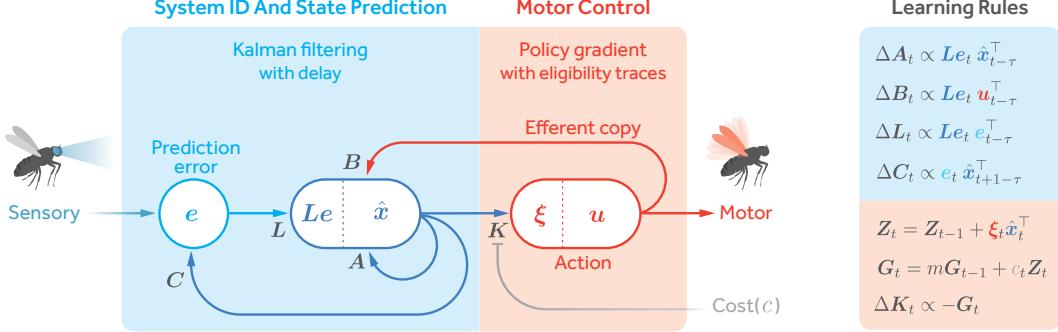| | [4] | [6] | [5] | [7] | [8] | Bio-OFC |
|---|---|---|---|---|---|---|
| delayed sensory feedback | ✗ | ✗ | ✗ | ✗ | ✗ | ✓ |
| control included | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| noise covariance agnostic | ✗ | ✓ | ✗ | ✗ | ✗ | ✓ |
| online system identification | ✗ | ✓ | ✗ | ✓ | ✓✗ | ✓ |
| local learning rules | N/A | ✓ | N/A | ✗ | ✓ | ✓ |
| tractable latent size | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ |
| absence of inner loop | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| single phase learning/execution | N/A | ✗ | N/A | ✓ | ✓ | ✓ |

Figure 1: **The circuit and learning rules of the Bio-OFC algorithm.** Our circuit is comprised of two main parts. First (in blue), the circuit performs Kalman filtering. Then (in red), the circuit performs control using policy gradients with eligibility traces. Triangular arrowheads denote synaptic connections and the flat arrowhead denotes the modulatory effect of the cost signal.

## 2 Background

We review classical Kalman estimation and control in this section, using boldface lowercase/uppercase letters for vectors/matrices and $\boldsymbol{I}$ for the identity matrix.

### 2.1 Problem formulation

We model the environment as a linear dynamical system driven by control input and perturbed by Gaussian noise. The true state of the system $\boldsymbol{x}$ is hidden and all the animal has access to are the observations $\boldsymbol{y}$ that are assumed to be linear functions of the state corrupted by Gaussian noise.

$$\text{dynamics:} \qquad \boldsymbol{x}_{t+1} = \boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t + \boldsymbol{v}_t \tag{1}$$

$$\text{observation:} \qquad \boldsymbol{y}_t = \boldsymbol{C}\boldsymbol{x}_t + \boldsymbol{w}_t \tag{2}$$

Here $\boldsymbol{v}_t \sim \mathcal{N}(0; \boldsymbol{V})$ and $\boldsymbol{w}_t \sim \mathcal{N}(0, \boldsymbol{W})$ are independent Gaussian random variables and the initial state has a Gaussian prior distribution $\boldsymbol{x}_0 \sim \mathcal{N}(\hat{\boldsymbol{x}}_0, \boldsymbol{\Sigma}_0)$.

The goal is to estimate the latent state $\hat{\boldsymbol{x}}$ in order to design a control $\boldsymbol{u}$ that minimizes expected cost

$$\text{expected cost:} \qquad J = \mathbb{E}\left[\sum_{t=0}^{T} c(\boldsymbol{x}_t, \boldsymbol{u}_t)\right] \tag{3}$$

$$\text{control:} \qquad \boldsymbol{u}_t = k(\hat{\boldsymbol{x}}_t) = \arg\min J \tag{4}$$

where $c(\boldsymbol{x}_t, \boldsymbol{u}_t)$ is the instantaneous cost associated with state $\boldsymbol{x}_t$ and action $\boldsymbol{u}_t$. As the environment dynamics is not known to the animal a priori, the parameters $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ must be learned online.

### 2.2 Kalman estimation and control

The Kalman filter [16, 2] is an estimator of the latent state $\hat{\boldsymbol{x}}_t$ (and its variance) via a weighted summation of the current observation and the prediction of the internal model based on prior measurements. However, in biologically realistic situations the sensory feedback is always delayed. This means that the control signal $\boldsymbol{u}_t$ has to be issued, and thus the state $\boldsymbol{x}_t$ estimated, before $\boldsymbol{y}_t$ has been observed. The appropriately modified Kalman filter computes the posterior probability distribution of $\boldsymbol{x}_t$ given observations $\boldsymbol{y}_{t-\tau}, ..., \boldsymbol{y}_0$ where $\tau \geq 1$ is the delay. We start here with the case of $\tau = 1$ for which the recursive updates of the mean $\hat{\boldsymbol{x}}_t$ and the variance $\boldsymbol{\Sigma}_t$ are

$$\hat{\boldsymbol{x}}_{t+1} = \boldsymbol{A}\hat{\boldsymbol{x}}_t + \boldsymbol{B}\boldsymbol{u}_t + \boldsymbol{L}_t(\boldsymbol{y}_t - \boldsymbol{C}\hat{\boldsymbol{x}}_t) \tag{5}$$

$$\boldsymbol{\Sigma}_{t+1} = (\boldsymbol{A} - \boldsymbol{L}_t\boldsymbol{C})\boldsymbol{\Sigma}_t\boldsymbol{A}^\top + \boldsymbol{V} \tag{6}$$

where $\boldsymbol{L}_t$ is known as the Kalman gain matrix which optimally combines the noisy observations $\boldsymbol{y}_t$ with the internal model and is given by

$$\boldsymbol{L}_t = \boldsymbol{A}\boldsymbol{\Sigma}_t\boldsymbol{C}^\top(\boldsymbol{C}\boldsymbol{\Sigma}_t\boldsymbol{C}^\top + \boldsymbol{W})^{-1}. \tag{7}$$

The Kalman filter is optimal in the sense that it minimizes the mean-squared error $\mathbb{E}[e_t \top e_t]$ with prediction error (innovation) $e_t = y_t - C\hat{x}_t$.

The output feedback law, also known as policy, (4) simplifies if the cost $J$ is quadratic in $u_t$ and $x_t$:

$$u_t = -K\hat{x}_t \tag{8}$$

and is known as linear-quadratic regulator (LQR). The control gain, $K$, is found by solving a matrix Riccati equation (cf. Supplementary Material). Linear policies have been successfully applied to a variety of control tasks [17].

## 3 Neural network representation for optimal feedback control

### 3.1 Inference

For now, let us assume that the system dynamics, the Kalman gain and the control gain $A, B, -C$, $L$, and $K$ are constant and known. Then, the latent state $\hat{x}$ can be obtained by the Kalman estimator Eq. (5). This algorithm naturally maps onto the network in Fig. 1 and Supplementary Fig. S1A with neural populations representing $\hat{x}$, $e := y - C\hat{x}$, and $u$ that are connected by synapses whose weights represent the elements of the matrices $A, B, -C$ and $L$. The computation of the control variable $u$ according to Eq. (S2) can be implemented by synapses whose weights represent the elements of the matrix $-K$.

If the sensory stimulus delay is $\tau = 1$, our network implements the Kalman prediction reviewed in the previous section. In case $\tau > 1$, the latent state must be recomputed throughout the delay period which requires a biologically implausible circuit (cf. Supplementary Fig. S1C). To overcome this, we adapt an alternative solution from online control [18, 19]. Specifically, we combine delayed measurements with the similarly delayed latent state estimation. Such inference can be performed by the network of the same architecture and adjusting the synaptic delay associated with matrix $C$ to match with the sensory delay, in accordance with the following expression:

$$\hat{x}_{t+1} = A\hat{x}_t + Bu_t + L\underbrace{(y_{t+1-\tau} - C\hat{x}_{t+1-\tau})}_{e_t} \tag{9}$$

As we show in Results (Fig. 3) this reduces predictive performance only modestly compared to the biologically unrealistic scheme. More details on the inference process and the temporal order in which our recurrent network performs the above steps is shown in Supplementary Fig. S1B.

### 3.2 Learning

#### 3.2.1 System identification and Kalman gain

Next, we turn our attention to the system identification/learning problem, which was not addressed by some of the previous proposals, cf. Table 1. Given a sequence of observations $\{y_0, \cdots, y_T\}$, we use a least squares approach [20] to find the parameters that minimize the mean-square prediction error $\frac{1}{T}\sum_{t=0}^{T} e_t^\top e_t$. We perform this optimization in an online manner, that is at each time-step $t+1$, after making the delayed observation $y_{t+1-\tau}$, we update the parameter estimates $\hat{A}, \hat{B}, \hat{C}$ using steps that minimize $e_t^\top e_t$, assuming that the state estimate and actions corresponding to prior observations (e.g. $\hat{x}_{t-\tau}, u_{t-\tau}$ etc.) are fixed. To obtain $L$ we would like to avoid solving the Riccati equation (6) as it requires matrix operations difficult to implement in biology. So, we use the same optimization procedure to update the Kalman gain $L$. Using Eq. (9) and explicitly writing out the matrix/vector indices as superscripts yields the following stochastic gradient with respect to $A$:

$$-\frac{\partial}{\partial A^{ij}} \frac{1}{2} \sum_k \left(e_t^k\right)^2 = -\sum_k e_t^k \frac{\partial e_t^k}{\partial A^{ij}} = -\sum_k e_t^k \frac{\partial \left(y_{t+1-\tau}^k - \sum_l C^{kl}\hat{x}_{t+1-\tau}^l\right)}{\partial A^{ij}} =$$

$$= \sum_{k,l} e_t^k C^{kl} \frac{\partial \left(\sum_m A^{lm}\hat{x}_{t-\tau}^m + \sum_n B^{ln}u_t^n + \sum_p L^{lp}e_t^p\right)}{\partial A^{ij}} =$$

$$= \sum_{k,l,m} e_t^k C^{kl}\delta^{li}\delta^{mj}\hat{x}_{t-\tau}^m = \sum_k e_t^k C^{ki}\hat{x}_{t-\tau}^j \tag{10}$$

4

Performing similar derivations for the other synaptic weights, our optimization procedure would rely on the following stochastic gradients:

$$-\nabla_{\boldsymbol{A}} \tfrac{1}{2} \boldsymbol{e}_t^\top \boldsymbol{e}_t = \boldsymbol{C}^\top \boldsymbol{e}_t \hat{\boldsymbol{x}}_{t-\tau}^\top \qquad\qquad -\nabla_{\boldsymbol{B}} \tfrac{1}{2} \boldsymbol{e}_t^\top \boldsymbol{e}_t = \boldsymbol{C}^\top \boldsymbol{e}_t \boldsymbol{u}_{t-\tau}^\top \qquad (11)$$

$$-\nabla_{\boldsymbol{L}} \tfrac{1}{2} \boldsymbol{e}_t^\top \boldsymbol{e}_t = \boldsymbol{C}^\top \boldsymbol{e}_t \boldsymbol{e}_{t-\tau}^\top \qquad\qquad -\nabla_{\boldsymbol{C}} \tfrac{1}{2} \boldsymbol{e}_t^\top \boldsymbol{e}_t = \boldsymbol{e}_t \hat{\boldsymbol{x}}_{t+1-\tau}^\top \qquad (12)$$

This yields a classical Hebbian rule between (a memory trace of) presynaptic activity $\hat{\boldsymbol{x}}_{t+1-\tau}$ and postsynaptic activity $\boldsymbol{e}_t$ for weights $\boldsymbol{C}$. However, it suggests non-local learning rules for $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{L}$, which runs contrary to biological requirements. We can circumvent this problem by replacing $\boldsymbol{C}^\top$ with $\boldsymbol{L}$, which corresponds to left-multiplication of the gradients with a positive definite matrix (see Supplementary Material Sec. D). This still decreases the mean-square prediction error under some mild initialization constraints on $\boldsymbol{C}$ and $\boldsymbol{L}$ and yields local plasticity rules, cf. Fig. 1,

$$\Delta \hat{\boldsymbol{A}}_t \propto \boldsymbol{L} \boldsymbol{e}_t\, \hat{\boldsymbol{x}}_{t-\tau}^\top \qquad (13) \qquad\qquad \Delta \hat{\boldsymbol{B}}_t \propto \boldsymbol{L} \boldsymbol{e}_t\, \boldsymbol{u}_{t-\tau}^\top \qquad (15)$$

$$\Delta \boldsymbol{L}_t \propto \boldsymbol{L} \boldsymbol{e}_t\, \boldsymbol{e}_{t-\tau}^\top \qquad (14) \qquad\qquad \Delta \hat{\boldsymbol{C}}_t \propto \boldsymbol{e}_t\, \hat{\boldsymbol{x}}_{t+1-\tau}^\top \qquad (16)$$

where the input current $\boldsymbol{L}\boldsymbol{e}_t$ is locally available at neurons representing $\hat{\boldsymbol{x}}$. The first three rules are local, but non-Hebbian, capturing correlations between presynaptic activity $\hat{\boldsymbol{x}}_{t-\tau}, \boldsymbol{u}_{t-\tau}, \boldsymbol{e}_{t-\tau}$ and postsynaptic current $\boldsymbol{L}\boldsymbol{e}_t$. Note that these updates do <u>not</u> require knowledge of the noise covariances $\boldsymbol{V}$ and $\boldsymbol{W}$, an advantage over previous work, cf. Table 1.

### 3.2.2  Control

Next, we consider optimal control a neural implementation of which was missing in most previous proposals, cf. Table 1. Traditionally, optimal control law is computed by iterating a matrix Riccati equation, cf. Supplementary Material, posing a difficult challenge for a biological neural implementation (but see [21]). To circumvent this problem, we propose to learn the controller weights $\boldsymbol{K}$ using a policy gradient method [12, 22] instead. Policy gradient methods directly parametrize a stochastic controller $\pi_{\boldsymbol{K}}(\boldsymbol{u}|\boldsymbol{x})$. Representing the total cost for a given trajectory $\tau$ as $c(\tau)$, they optimize the parameters $\boldsymbol{K}$ by performing gradient descent on the expected cost $J = \mathbb{E}_{\pi_{\boldsymbol{K}}}[c(\tau)] = \mathbb{E}_{\pi_{\boldsymbol{K}}}\left[\sum_{t=0}^{T} c_t\right]$.

$$\nabla_{\boldsymbol{K}} J = \int c(\tau) \nabla_{\boldsymbol{K}} \pi_{\boldsymbol{K}}(\tau) d\tau = \mathbb{E}_{\pi_{\boldsymbol{K}}}\left[c(\tau) \nabla_{\boldsymbol{K}} \log \pi_{\boldsymbol{K}}(\tau)\right] \qquad (17)$$

$$= \mathbb{E}_{\pi_{\boldsymbol{K}}}\left[\left(\sum_{t=0}^{T} c_t\right)\left(\sum_{s=0}^{T} \nabla_{\boldsymbol{K}} \log \pi_{\boldsymbol{K}}(\boldsymbol{u}_s|\boldsymbol{x}_s)\right)\right] \qquad (18)$$

The term in square brackets is an unbiased estimator of the gradient and can be used to perform stochastic gradient decent. As already hinted at by [12], due to causality, costs $c_t$ are not affected by later controls $\boldsymbol{u}_s$, $s > t$, and the variance of the estimator can be reduced by excluding those terms, which yields

$$\Delta \boldsymbol{K} \propto -\sum_{t=0}^{T} c_t \left(\sum_{s=0}^{t} \nabla_{\boldsymbol{K}} \log \pi_{\boldsymbol{K}}(\boldsymbol{u}_s|\boldsymbol{x}_s)\right). \qquad (19)$$

We simply keep an eligibility trace of the past, $\boldsymbol{Z}_t = \sum_{s=0}^{t} \nabla_{\boldsymbol{K}} \log \pi_{\boldsymbol{K}}(\boldsymbol{u}_s|\boldsymbol{x}_s)$, and perform parameter updates $\Delta \boldsymbol{K} \propto -c_t \boldsymbol{Z}_t$ at each time step $t$. A similar update rule has been suggested by [23, 24] for the infinite horizon case. Global convergence of policy gradient methods for linear quadratic regulator has been recently studied by Fazel *et al.* [25].

We assume the output of neurons encoding control $\boldsymbol{u}$ is perturbed by Gaussian noise

$$\boldsymbol{u}_t = -\boldsymbol{K}\hat{\boldsymbol{x}}_t - \boldsymbol{\xi}_t \quad \text{with} \quad \boldsymbol{\xi}_t \sim \mathcal{N}(0, \sigma^2 \boldsymbol{I}). \qquad (20)$$

The synapses are endowed with a synaptic tag [26] $\boldsymbol{Z}$, an eligibility trace that tracks correlations between pre-synaptic activity $\hat{\boldsymbol{x}}$ and post-synaptic noise $\boldsymbol{\xi}$. It is reset to zero at the beginning of each trajectory, though instead of a hard reset it could also softly decay with a time constant of the same order $\mathcal{O}(T)$ as trajectory duration [27]. The weight update assigns cost $c_t$ to the synapses according to their eligibility $\boldsymbol{Z}_t$. The cost is e.g. provided by a diffuse neuromodulatory signal such as dopamine.

The optional use of momentum $m \in [0, 1)$ adds a low-pass filter to the synaptic plasticity cascade:

$$\boldsymbol{Z}_t = \boldsymbol{Z}_{t-1} + \boldsymbol{\xi}_t \hat{\boldsymbol{x}}_t^\top \quad \left( = \sigma^2 \sum_{s=0}^t \nabla_{\boldsymbol{K}} \log \pi_{\boldsymbol{K}}(\boldsymbol{u}_s | \hat{\boldsymbol{x}}_s) \right) \tag{21}$$

$$\boldsymbol{G}_t = m \boldsymbol{G}_{t-1} + c_t \boldsymbol{Z}_t \tag{22}$$

$$\Delta \boldsymbol{K}_t \propto -\boldsymbol{G}_t \tag{23}$$

## 4   Experiments

In this section we look at three different experiments to demonstrate various features of the Bio-OFC algorithm. In Sec. 4.1, we look at a discrete-time double integrator and discuss how our approach performs not only Kalman filtering (learning the optimal Kalman gain), but also full system-ID in the open-loop setting (system-ID followed by control) as well as in the more challenging closed-loop setting (simultaneous system-ID and control). In each case, we provide quantitative comparisons and discuss the effect of increased delay. In Secs. 4.2 and 4.3, we apply our methodology to two biologically relevant control tasks, that of reaching movements and flight.

### 4.1   Discrete-time double integrator

As a simple test case, we follow along the lines of [28] and consider the classic problem of a discrete-time double integrator with the dynamical model

$$\boldsymbol{x}_{t+1} \sim \mathcal{N}(\boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t, \boldsymbol{V}) \quad \text{where} \quad \boldsymbol{A} = \begin{pmatrix} 1 & 1 \\ 0 & 1 \end{pmatrix}, \quad \boldsymbol{B} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \quad \boldsymbol{V} = \begin{pmatrix} .01 & 0 \\ 0 & .01 \end{pmatrix}. \tag{24}$$

Such a system models the position and velocity (respectively the first and second components of the state) of a unit mass object under force $u$. As an instance of LQR, we can try to steer this system to reach point $(0,0)^\top$ from initial condition $\boldsymbol{x}_0 = (-1, 0)^\top$ without expending much force:

$$J = \sum_{t=0}^T \boldsymbol{x}_t^\top \boldsymbol{Q} \boldsymbol{x}_t + R \sum_{t=0}^{T-1} u_t^2 \quad \text{where} \quad \boldsymbol{Q} = \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, \quad R = 1, \quad T = 10 \tag{25}$$

We assume that the initial state estimate $\hat{\boldsymbol{x}}_0$ is at the true initial state ($\hat{\boldsymbol{x}}_0 = \hat{\boldsymbol{C}}^+ \boldsymbol{C} \boldsymbol{x}_0$).

We go beyond LQR (cf. Supplementary Fig. S2A) and assume $\boldsymbol{x}$ is not directly observable (cf. Supplementary Fig. S2B,C), but we merely have access to noisy observations, $\boldsymbol{y} \sim \mathcal{N}(\boldsymbol{C}\boldsymbol{x}, \boldsymbol{W})$. We consider two observation models, one where the state is only observed with some uncorrelated noise, and one where the observation noise covariance is not diagonal and an additional mixture of the 2 state components is observed:

$$\text{LDS1:} \ \boldsymbol{C} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}, \boldsymbol{W} = \begin{pmatrix} .04 & 0 \\ 0 & .25 \end{pmatrix} \qquad \text{LDS2:} \ \boldsymbol{C} = \begin{pmatrix} 1 & 0 \\ 0 & -1 \\ .5 & .5 \end{pmatrix}, \boldsymbol{W} = \begin{pmatrix} .04 & .09 & 0 \\ .09 & .25 & 0 \\ 0 & 0 & .04 \end{pmatrix} \tag{26}$$

We denote these two systems as linear dynamical systems 1 and 2 (LDS1 and LDS2).

**Learning the Kalman gain.**   A major advantage of our work is that it does not require knowledge of the covariance matrices $\boldsymbol{V}$ and $\boldsymbol{W}$ to determine the Kalman filter gain $\boldsymbol{L}$. We studied this using LDS1 in a scenario where the observation noise varies; changing the covariance matrix $\boldsymbol{W}$ from $\mathrm{diag}(.04, .25)$ to $\mathrm{diag}(.04, .01)$ after 2500, to $\mathrm{diag}(.01, .01)$ after 5000, and back to $\mathrm{diag}(.04, .25)$ after 7500 episodes. In this experiment we fix $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ at the ground truth, and initialize $\boldsymbol{L}$ at the optimal value for $\boldsymbol{W} = \mathrm{diag}(.04, .25)$, and updated the latter according to Eq. (14). Fig. 2 shows how the elements of the filter matrix $\boldsymbol{L}$ adapt in time, to optimize performance as measured by the mean squared prediction error. The learning rate was tuned to minimize the average MSE over all episodes. Although the Kalman gain can be slow to converge in some cases (Fig. 2A), performance is quickly close to optimal (Fig. 2B).

Fig. 3 shows the achievable optimal control cost as a function of delay for four different controllers: the optimal linear-quadratic-Gaussian (LQG) controller that uses time-dependent gains, the biologically implausible ANN (cf. Supplementary Fig. S1C) that uses time-invariant gains, a model-free[1] approach

---

[1]Note that LQG uses the model for state inference as well as for control. Our Bio-OFC uses the model only to infer the latent state, but uses a model-free controller. In contrast, model-based reinforcement learning algorithms typically assume knowledge of the current state and use the model only for control.
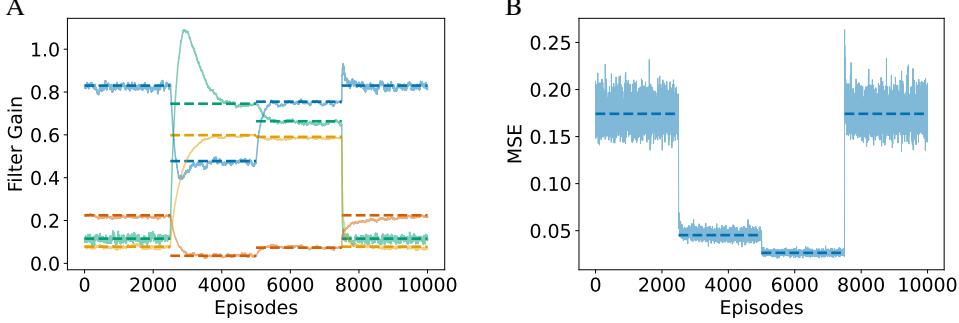
Figure 2: **Bio-OFC adapts to changing noise statistics.** **(A)** Filter gain (colors denote different elements of the gain matrix) and **(B)** mean squared prediction error (MSE) in the simple LQG task with 2 latent dimensions and 2-d observations (LDS1, see text for details). Solid lines show the mean over 20 runs. Dashed lines indicate the optimal filter gain and the corresponding average MSE. After 2500 episodes the observation noise covariance $\boldsymbol{W}$ decreased to $\mathrm{diag}(.04, .01)$, after 5000 to $\mathrm{diag}(.01, .01)$, and after 7500 back to $\mathrm{diag}(.04, .25)$.



Figure 3: **Optimal cost of Bio-OFC is close to that of LQG for various delays.** Optimal cost as function of measurement delay obtained by different methods. We used the simple LQG tasks with 2 latent dimensions and **(A)** 2-d observations (LDS1, see text for details) and **(B)** 3-d observations (LDS2). LQG uses the optimal time-dependent filter and feedback gains, ANN uses the network of Supplementary Fig. S1C, Bio-OFC the network of Fig. 1 (cf. Supplementary Figs. S1B and S2D), and model-free directly maps from noisy delayed observations to control, $\boldsymbol{u}_t = \boldsymbol{K}\boldsymbol{y}_{t-\tau}$ (cf. Supplementary Fig. S2E). Shown is the mean cost $\pm$ SEM over 10000 episodes.

using policy gradient method applied directly to observations (inset), and Bio-OFC that update the current state estimate $\hat{\boldsymbol{x}}_t$ directly using the delayed measurement $\boldsymbol{y}_{t-\tau}$ (cf. Eq. (9) and Fig. 1). For the sake of biological plausibility, we imposed time-invariant gains and direct state estimate updates based on delayed measurement. These results clearly demonstrate that Bio-OFC is robust to sensory delays. Specifically, while it is expected that Bio-OFC will not learn a solution quite as good as LQG, our results show that the solution found by it is not very far off. In general, a model can be useful in two ways: It facilitates a filtered estimate that is useful even in the absence of measurement delays (see LQG vs model-free for delay 0 in Fig. 3), and in the presence of delays the model helps to bridge the gap by predicting forward in time.

**Full system identification.** We next considered the case of system identification, i.e. learning the weight matrices $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ in addition to $\boldsymbol{L}$, using Eqs. (13-16). We initialized $\boldsymbol{A}$ and $\boldsymbol{B}$ with small random numbers drawn from $\mathcal{N}(0, 0.01)$. For LDS1, $\boldsymbol{C}$ and $\boldsymbol{L}$ were initialized as diagonal dominated random matrices, with diagonal elements drawn uniformly randomly from $[0.5, 1]$ and off-diagonal ones from $[0, 0.5]$. For LDS2, they were drawn from $\mathcal{N}(0, 0.01)$ under the constraint that the symmetric part of $\boldsymbol{LC}$ has positive eigenvalues. Controls $u_t$ were drawn from $\mathcal{N}(0, 0.25)$. Fig. 4A,B show how the mean squared prediction error converges to the optimal values from Fig. 3 (dashed horizontal lines). We also considered a Bio-OFC that learns in environment LDS2, but over-represents the latent space, assuming it is 3-d instead of 2-d. The square matrices $\boldsymbol{C}$ and $\boldsymbol{L}$

were initialized analogously to LDS1. Fig. 4C shows that this over-representation does not affect performance.
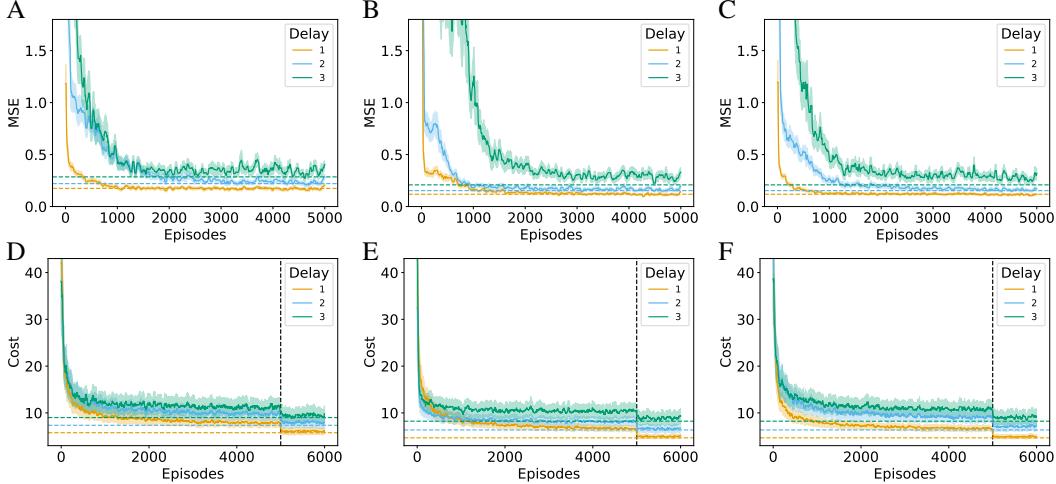


Figure 4: **Open-loop Bio-OFC converges to optimally achievable MSE and cost given the delay.** The solid (shaded) curve depicts the mean ($\pm$SEM) of 20 runs with different random initial weights, smoothed using a running median with window size 51. The dashed horizontal lines show the asymptotic values from Fig. 3. The dashed vertical line depicts the time when learning was stopped and the exploratory stochasticity in the controller removed. Mean squared prediction error during system identification, using plasticity rules Eqs. (13-16) for the filter and random Gaussian controller input, as function of episodes for **(A)** 2-d observations (LDS1), **(B)** 3-d observations (LDS2) and **(C)** 3-d observations (LDS2) with an over-representing Bio-OFC that assumes 3 instead of the actual 2 latent dimensions. Cost as function of episodes, using plasticity rules Eqs. (21-23) for the controller, for **(D)** 2-d observations (LDS1, see text for details), **(E)** 3-d observations (LDS2) and **(F)** 3-d observations (LDS2) with an over-representing Bio-OFC that assumes 3 instead of the actual 2 latent dimensions.

**System identification followed by control.** After performing system identification for 5000 episodes, we kept $A, B, C, L$ fixed and transitioned to learning the controller $K$ using Eqs. (21-23) for another 5000 episodes with controller noise $\xi \sim \mathcal{N}(0, 0.04)$, cf. Eq. (20) and Fig. 4D-F. The stochasticity in the controller results in an excess cost. Using a deterministic controller ($\xi = 0$) for another 1000 episodes reveals that the network converged to the optimal cost from Fig. 3 (dashed horizontal lines). We used two learning rates, one for $A, B, C, L$ and one for $K$ with momentum $m = 0.99$ for the latter. Learning rates that quickly yield good final performance were chosen by minimizing the sum of average reward during and after learning using Optuna [29], a hyperparameter optimization framework freely available under the MIT license. Different noise levels $\sigma$ and momenta $m$ are considered in Supplementary Figs. S6 and S7 respectively.

**Simultaneous system ID and control.** Performing system identification prior to learning a controller is known as open-loop adaptive control and a common approach in control theory. Recent advances in the control community tackle the more challenging problem of closed-loop control [14], [2] simultaneously identifying both while controlling the system. Our network is capable of closed-loop control with no separate phases for system identification and control optimization necessary. In contrast to our work, the only other proposed neural implementation that also includes control [6] requires separate phases for system-ID and control. Fig. 5 shows how the control cost evolves in time when the weights are updated according to Eqs. (13-16) and Eqs. (21-23) while using the controller designed inputs of Eq. (20). Again, using a deterministic controller ($\xi = 0$) for another 1000 episodes reveals that the network converged to the optimal cost. Different noise levels $\sigma$ and momenta $m$ are considered in Supplementary Figs. S8 and S9 respectively.

---

[2] When a controller designs the inputs based on the history of inputs and observations, the inputs become highly correlated with the past process noise sequences, which prevents consistent and reliable parameter estimation with standard system identification techniques.
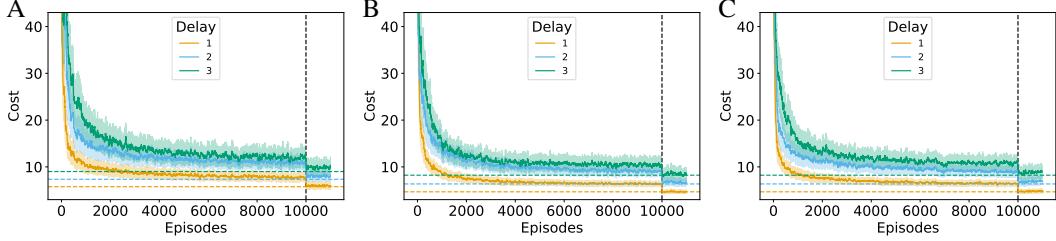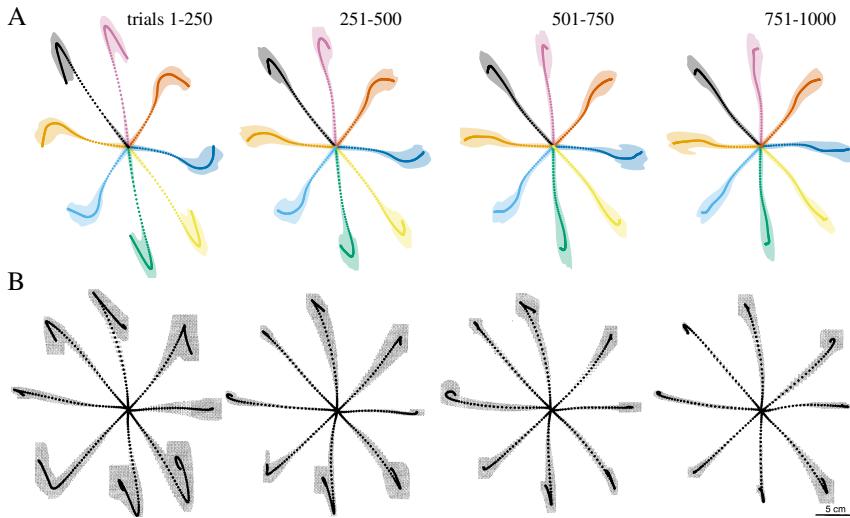
Figure 5: **Closed-loop Bio-OFC converges to optimally achievable cost given the delay.** The solid (shaded) curve depicts the mean (±SEM) of 20 runs with different random initial weights, smoothed using a running median with window size 51. The dashed horizontal lines show the asymptotic values from Fig. 3. The dashed vertical line depicts the time when learning was stopped and the exploratory stochasticity in the controller removed. Cost as function of episodes, using plasticity rules Eqs. (13-16) for the filter and Eqs. (21-23) for the controller simultaneously, for **(A)** 2-d observations (LDS1, see text for details), **(B)** 3-d observations (LDS2) and **(C)** 3-d observations (LDS2) with an over-representing Bio-OFC that assumes 3 instead of the actual 2 latent dimensions.



Figure 6: **Bio-OFC captures human performance when learning to adapt to a force field. (A)** Model trajectories during training. Performance plotted during the first, second, third, and final 250 targets. Dots show the mean and are 10 ms apart, shaded area shows a kernel density estimate thresholded at 0.04. **(B)** Averages±SD of human hand trajectories during training [1]. Copyright ©1994 Society for Neuroscience.

## 4.2 Reaching movements

To connect back to a biological sensory-motor control task, we considered the task of making reaching movements in the presence of externally imposed forces from a mechanical environment [1] (Supplementary Material). We initialized the weights of our network to the values that are optimal in a null force field, using a unit time of 10 ms and a sensory delay of 50 ms (i.e. $\tau = 5$), as has been measured experimentally [30]. Bio-OFC successfully captures the characteristics of human trajectories in the null field as well as the force field, cf. Fig. S10. Bio-OFC adapts to the force field by updating its weights according to plasticity rules Eqs. (13-16) for the filter and Eqs. (21-23) for the controller. Figs. 6 and S11 show that this captures human performance during the training period. Switching off learning in the controller yields virtually identical results, cf. Fig. S12, thus learning is driven primarily by changes in the estimator. Using signal-dependent motor noise in the plant [31], which increases with the magnitude of the control signal, also yields similar results, cf. Figs. S13-S14.

9

### 4.3 Simplified winged flight

For our final example, we designed an OpenAI gym [32] environment which simulates winged flight in 2-d with simplified dynamics (cf. Fig. S15). Here, the agent controls the flapping frequency of each wing individually, producing an impulse up and away from the wing (i.e. direction up and left when flapping the right wing). The agent receives sensory stimuli which are delayed by $100\,\mathrm{ms}$ (equivalent to $\tau = 5$ time-steps of the simulation). The goal of the agent is to fly to a fixed target and stabilize itself against gravity, the environment wind, and stochastic noise in the control system. The agent suffers a cost that is proportional to the L1 distance to the target, and the L1 magnitude of the control variables. This L1 cost was chosen to verify the flexibility of our algorithm when the cost deviates from the assumptions of LQR, where the cost is quadratic. We compare the performance of Bio-OFC to policy gradient. We find that, because of the delay, the agent trained with policy gradient overshoots the target and needs to backtrack, cf. Fig. S16a. However, the agent trained with Bio-OFC flies directly towards the target with no significant overshoot, cf. Fig. S16b. For more details and a video demonstration (`gym-fly-demo.mp4`) see the Supplementary Material.

## 5   Discussion

In this work, we developed a biologically plausible neural algorithm for sensory-motor control that learns a model of the world, and uses the ability to preview future, through a neurally plausible version of the Kalman filter, to learn an appropriate control policy. This neural circuit has the capacity to build an adequate representation of the appropriate state space, can deal with sensory delays, and actively explores the action space to execute appropriate control strategies.

We used a model-free controller, primarily due to its simplicity and biological plausibility. A model-based controller would need access to the model, i.e. weight matrices $A$ and $B$, which can result in a weight transport problem. However, model-based control has advantages such as higher sample efficiency and the ability to be transferable to other goals and tasks. An interesting question for future work would be how to combine state estimation via the Kalman filter with model-based control in a biologically plausible manner.

One limitation of this work is that it is in the framework of linear control theory. Locally linearized dynamics [33] has been suggested to generalize the Kalman filter. The inputs could also be processed using additional neural network layers to obtain a representation that renders the dynamics linear [34]. In several normative approaches towards neurally plausible representation learning, simply constraining neural activity to be nonnegative while retaining the same objective functions, allowed one to move from, say, PCA [35] to clusters [36] and manifolds [37]. Our work could be the starting point of a similar generalization.

We considered a uniform delay for all stimuli. In the case of motor control, proprioceptive feedback is faster than visual feedback [30]. Our model readily extends to the case of various, but known, delays for different modalities. The prediction error coding neurons merely need to combine predictions and measurements adequately, i.e. the synaptic delay associated with prediction has to match with the sensory delay for that modality. We believe learning the appropriate delay could be implemented by extending the state space using lag vectors [2], which we leave for future work.

In line with overall brain architecture [38] and the predictive coding framework [10], our model suggests the brain employs a recurrent network to generate predictions and actions by constantly attempting to match incoming sensory inputs with top-down predictions [11]. Our model can also be mapped to brain regions putatively contributing to optimal feedback control [30, 39]. Specifically, it has been proposed that the cerebellum performs system identification, parietal cortex performs state estimation, and primary and premotor cortices implement the optimal control policy by transforming state estimates into motor commands. Also, basal ganglia may be involved in processing cost/reward [39, 40].

This work proposed a concrete neural architecture that takes up the challenge of online control in a changing world [41], with delay, and using biologically plausible synaptic update rules. The learning algorithm performs well for several tasks, even with some drastic approximations. Future exploration of its limitations would provide further insights into the general nature of biologically constrained control systems.

## Acknowledgments and Disclosure of Funding

## References

[1] R. Shadmehr and F. A. Mussa-Ivaldi. Adaptive representation of dynamics during learning of a motor task. Journal of Neuroscience, 14(5):3208–3224, 1994.

[2] E. Todorov and M. I. Jordan. Optimal feedback control as a theory of motor coordination. Nature Neuroscience, 5(11):1226–1235, 2002.

[3] S. H. Scott. Optimal feedback control and the neural basis of volitional motor control. Nature Reviews Neuroscience, 5(7):532–545, 2004.

[4] S. Deneve, J.-R. Duhamel, and A. Pouget. Optimal sensorimotor integration in recurrent cortical networks: a neural implementation of Kalman filters. Journal of Neuroscience, 27(21):5744–5756, 2007.

[5] R. Wilson and L. Finkel. A neural implementation of the Kalman filter. In Advances in Neural Information Processing Systems, volume 22, pages 2062–2070, 2009.

[6] R. Linsker. Neural network learning of optimal Kalman prediction and control. Neural Networks, 21(9):1328–1343, 2008.

[7] A. Kutschireiter, S. C. Surace, H. Sprekeler, and J.-P. Pfister. Nonlinear Bayesian filtering and learning: a neuronal dynamics for perception. Scientific Reports, 7(1):1–13, 2017.

[8] B. Millidge, A. Tschantz, A. Seth, and C. Buckley. Neural Kalman filtering. arXiv:2102.10021, 2021.

[9] D. McNamee and D. M. Wolpert. Internal models in biological control. Annual Review of Control, Robotics, and Autonomous Systems, 2:339–364, 2019.

[10] R. P. Rao and D. H. Ballard. Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. Nature Neuroscience, 2(1):79–87, 1999.

[11] A. Clark. Whatever next? predictive brains, situated agents, and the future of cognitive science. Behavioral and Brain Sciences, 36(3):181–204, 2013.

[12] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. Machine Learning, 8(3-4):229–256, 1992.

[13] H. S. Seung. Learning in spiking neural networks by reinforcement of stochastic synaptic transmission. Neuron, 40(6):1063–1073, 2003.

[14] S. Lale, K. Azizzadenesheli, B. Hassibi, and A. Anandkumar. Logarithmic regret bound in partially observable linear dynamical systems. In Advances in Neural Information Processing Systems, volume 33, pages 20876–20888, 2020.

[15] B. Lee and A. Lamperski. Non-asymptotic closed-loop system identification using autoregressive processes and Hankel model reduction. In 2020 59th IEEE Conference on Decision and Control (CDC), pages 3419–3424, 2020.

[16] R. E. Kalman. A new approach to linear filtering and prediction problems. Journal of Basic Engineering, 82(1):35–45, 1960.

[17] A. Rajeswaran, K. Lowrey, E. V. Todorov, and S. M. Kakade. Towards generalization and simplicity in continuous control. In Advances in Neural Information Processing Systems, volume 30, 2017.

[18] H. L. Alexander. State estimation for distributed systems with sensing delay. In Data Structures and Target Classification, volume 1470, pages 103–111. International Society for Optics and Photonics, 1991.

[19] T. D. Larsen, N. A. Andersen, O. Ravn, and N. K. Poulsen. Incorporation of time delayed measurements in a discrete-time Kalman filter. In Proceedings of the 37th IEEE Conference on Decision and Control, volume 4, pages 3972–3977. IEEE, 1998.

[20] M. Verhaegen and V. Verdult. Filtering and system identification: a least squares approach. Cambridge University Press, 2007.

[21] J. Friedrich and M. Lengyel. Goal-directed decision making with spiking neurons. Journal of Neuroscience, 36(5):1529–1546, 2016.

[22] R. S. Sutton, D. A. McAllester, S. P. Singh, Y. Mansour, et al. Policy gradient methods for reinforcement learning with function approximation. In Advances in Neural Information Processing Systems, volume 99, pages 1057–1063, 1999.

[23] H. Kimura, M. Yamamura, and S. Kobayashi. Reinforcement learning by stochastic hill climbing on discounted reward. In Machine Learning Proceedings 1995, pages 295–303. Elsevier, 1995.

[24] J. Baxter and P. L. Bartlett. Infinite-horizon policy-gradient estimation. Journal of Artificial Intelligence Research, 15:319–350, 2001.

[25] M. Fazel, R. Ge, S. Kakade, and M. Mesbahi. Global convergence of policy gradient methods for the linear quadratic regulator. In International Conference on Machine Learning, pages 1467–1476. PMLR, 2018.

[26] R. L. Redondo and R. G. Morris. Making memories last: the synaptic tagging and capture hypothesis. Nature Reviews Neuroscience, 12(1):17–30, 2011.

[27] J. Friedrich, R. Urbanczik, and W. Senn. Spatio-temporal credit assignment in neuronal population learning. PLoS Computational Biology, 7(6):e1002092, 2011.

[28] B. Recht. A tour of reinforcement learning: The view from continuous control. Annual Review of Control, Robotics, and Autonomous Systems, 2:253–279, 2019.

[29] T. Akiba, S. Sano, T. Yanase, T. Ohta, and M. Koyama. Optuna: A next-generation hyperparameter optimization framework. In Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2019.

[30] S. H. Scott. The computational and neural basis of voluntary motor control and planning. Trends in Cognitive Sciences, 16(11):541–549, 2012.

[31] C. M. Harris and D. M. Wolpert. Signal-dependent noise determines motor planning. Nature, 394(6695):780–784, 1998.

[32] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. OpenAI Gym, 2016.

[33] S. J. Julier and J. K. Uhlmann. New extension of the Kalman filter to nonlinear systems. In Signal processing, sensor fusion, and target recognition VI, volume 3068, pages 182–193. International Society for Optics and Photonics, 1997.

[34] B. Lusch, J. N. Kutz, and S. L. Brunton. Deep learning for universal linear embeddings of nonlinear dynamics. Nature Communications, 9(1):1–10, 2018.

[35] C. Pehlevan, T. Hu, and D. B. Chklovskii. A Hebbian/anti-Hebbian neural network for linear subspace learning: A derivation from multidimensional scaling of streaming data. Neural Computation, 27(7):1461–1495, 2015.

[36] C. Pehlevan and D. B. Chklovskii. A normative theory of adaptive dimensionality reduction in neural networks. In Advances in Neural Information Processing Systems, volume 28, pages 2269–2277, 2015.

[37] A. M. Sengupta, C. Pehlevan, M. Tepper, A. Genkin, and D. B. Chklovskii. Manifold-tiling localized receptive fields are optimal in similarity-preserving neural networks. In Advances in Neural Information Processing Systems, volume 31, 2018.

[38] H. Sohn, N. Meirhaeghe, R. Rajalingham, and M. Jazayeri. A network perspective on sensorimotor learning. Trends in Neurosciences, 2020.

[39] R. Shadmehr and J. W. Krakauer. A computational neuroanatomy for motor control. Experimental Brain Research, 185(3):359–381, 2008.

[40] D. I. Todorov, R. A. Capps, W. H. Barnett, E. M. Latash, T. Kim, K. C. Hamade, S. N. Markin, I. A. Rybak, and Y. I. Molkov. The interplay between cerebellum and basal ganglia in motor adaptation: A modeling study. PLoS ONE, 14(4):e0214926, 2019.

[41] N. Agarwal, E. Hazan, A. Majumdar, and K. Singh. A regret minimization approach to iterative learning control. In International Conference on Machine Learning, volume 139, pages 100–109. PMLR, 2021.

[42] P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson, K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey, İ. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa, P. van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. Nature Methods, 17:261–272, 2020.

[43] E. Catto. Box2d: A 2d physics engine for games. https://box2d.org/, 2011.

# Supplementary material

## A    Linear-quadratic regulator (LQR)

The optimal control problem is to determine an output feedback law that minimizes the expected value of a cost criterion. If the cost $J$ is quadratic, the optimal output feedback is a linear control law known as linear-quadratic regulator (LQR).

$$J = \boldsymbol{x}_T^\top \boldsymbol{Q} \boldsymbol{x}_T + \sum_{t=0}^{T-1} \left( \boldsymbol{x}_t^\top \boldsymbol{Q} \boldsymbol{x}_t + \boldsymbol{u}_t^\top \boldsymbol{R} \boldsymbol{u}_t \right) \tag{S1}$$

$$\boldsymbol{u}_t = -\boldsymbol{K}_t \boldsymbol{x}_t \quad \text{with control gain} \quad \boldsymbol{K}_t = (\boldsymbol{B}^\top \boldsymbol{P}_{t+1} \boldsymbol{B} + \boldsymbol{R})^{-1} \boldsymbol{B}^\top \boldsymbol{P}_{t+1} \boldsymbol{A} \tag{S2}$$

where $\boldsymbol{P}_t$ is determined by the dynamic Riccati equation that runs backwards in time

$$\boldsymbol{P}_{t-1} = \boldsymbol{A}^\top \boldsymbol{P}_t \boldsymbol{A} - (\boldsymbol{A}^\top \boldsymbol{P}_t \boldsymbol{B}) \left( \boldsymbol{R} + \boldsymbol{B}^\top \boldsymbol{P}_t \boldsymbol{B} \right)^{-1} (\boldsymbol{B}^\top \boldsymbol{P}_t \boldsymbol{A}) + \boldsymbol{Q} \tag{S3}$$

from terminal condition $\boldsymbol{P}_T = \boldsymbol{Q}$. Linear-quadratic-Gaussian (LQG) control uses the Kalman estimate $\hat{\boldsymbol{x}}$ in the controller, $\boldsymbol{u}_t = -\boldsymbol{K}_t \hat{\boldsymbol{x}}_t$.

## B    Experimental details

The experiments to produce the figures of the paper were performed on a Linux-based (CentOS) desktop with Intel Xeon CPU E5-2643 v4 @ 3.40GHz (6 cores) and 128 GB of RAM. No usage of a GPU was made. We used SciPy's [42] minimize function (with the default BFGS algorithm) to optimize the learning rate for Fig. 2 (which took 98 s) and the optimal gains $\boldsymbol{K}$ and $\boldsymbol{L}$ for Fig. 3 (which, dependent on the delay, took from few seconds up to 2 minutes).

To produce Figs. 4 and 5 (also supporting Figs. S6-S9), 20 parallel runs with different random seeds for 10000+1000 episodes took 9.3 s for open loop and 14.0 s for closed loop control, largely irrespective of the considered LDS and delay. The learning rates used in those experiments were obtained earlier with Optuna [29]. This hyperparameter optimization was performed on a linux-based (CentOS) cluster with Intel Xeon CPU E5-2680 v4 @ 2.40GHz (14 cores) and 512 GB of RAM, dedicating an individual node to each combinatorial choice of LDS, control setting (open/closed loop), and delay. Using a computing budget of 1000 trials optimization took about 130 min for open loop and 190 min for closed loop control.

Computational complexity of our algorithms are the same as that of policy gradient and Kalman filtering. That is, our approximations and biological implementation do not change the computational complexity of these methods.

## C    Code

Code to reproduce the figures in the paper can be found in the GitHub repository `https://github.com/j-friedrich/neuralOFC`.
Requirements: `python`, `matplotlib`, `numpy`, `scipy`
The hyperparameters obtained with optuna [29] are provided in the subdirectory `results`.
To recreate a figure run the corresponding script. Figures will be saved in the subdirectory `fig`.

## D    Learning rules

The marginal mean-squared error decreases if the angle between the gradient $\boldsymbol{g}_\theta$ for parameter $\boldsymbol{\theta}$ and the update $\Delta\boldsymbol{\theta}$ is less than $90°$, i.e. $\boldsymbol{g}_\theta^\top \Delta\boldsymbol{\theta} > 0$. Thus to obtain $\Delta\boldsymbol{\theta}$ we can left-multiply the gradient $\boldsymbol{g}_\theta$ by any real square matrix $\boldsymbol{M}$ that is positive definite, i.e. its symmetric part $\frac{1}{2}(\boldsymbol{M} + \boldsymbol{M}^\top)$ has positive real eigenvalues, thus that $\boldsymbol{g}_\theta^\top (\boldsymbol{M} + \boldsymbol{M}^\top) \boldsymbol{g}_\theta > 0$ while the anti-symmetric part always satisfies $\boldsymbol{g}_\theta^\top (\boldsymbol{M} - \boldsymbol{M}^\top) \boldsymbol{g}_\theta = 0$. Replacing $\boldsymbol{C}^\top$ with $\boldsymbol{L}$ to obtain the learning rules, Eqs. (13-14), corresponds to multiplication of the gradients, Eqs. (11,12), by the product of $\boldsymbol{L}$ and the Moore-Penrose inverse of $\boldsymbol{C}^\top$, $\boldsymbol{M} = \boldsymbol{L}\boldsymbol{C}^{\top+}$. We therefore initialize $\boldsymbol{C}$ and $\boldsymbol{L}$ in such a way that $\boldsymbol{L}\boldsymbol{C}^{\top+}$ is positive definite (but not necessarily symmetric).
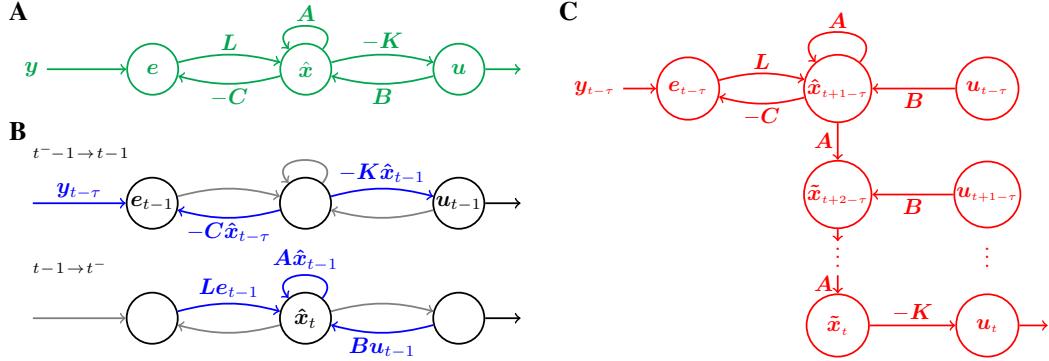
Figure S1: **Schematic of our proposed neural networks for OFC and alternative artificial neural network.** (A) Schematic of our proposed neural network (Bio-OFC). Nodes are annotated with the quantity they represent in their firing rates; edges are annotated with synaptic weights. (B) Input currents when updating $e_{t-1}$, $u_{t-1}$ (top), and $\hat{x}_t$ (bottom). $\hat{x}_t$ is updated directly using the delayed measurement $y_{t-\tau}$, even for delays $\tau > 1$. (C) The alternative artificial neural network (ANN) that updates the past estimate $\hat{x}_{t+1-\tau}$, and predicts forward in time to estimate $\tilde{x}_{t+2-\tau}, ..., \tilde{x}_t$, would require biologically implausible weight copies of $A$ and $B$, as well as a memory of past controls $u_{t-\tau}, ..., u_{t-1}$.



Figure S2: **Probabilistic graphical models of MDPs and POMDPs.** Grey colored nodes are observed, black and blue colored edges depend on the environment and actor, respectively. (A) MDP; the state is observed directly and the optimal action $u_t$ depends only on the current state $x_t$. (B) POMDP; the state is only partially observable and the action $u_t$ depends on the history of observations $y_0, ..., y_t$. (C) POMDP with delayed observations (here $\tau = 1$); the action $u_t$ depends on the limited history of available observations $y_0, ..., y_{t-\tau}$. (D) Model-based controller that produces action $u_t$ based on the current state $\hat{x}_t$ of the internal model that effectively summarizes past observations. Our Bio-OFC is an instance of such a controller. (E) Model-free memory-less controller that produces action $u_t$ based on the currently available delayed observation $y_{t-\tau}$.
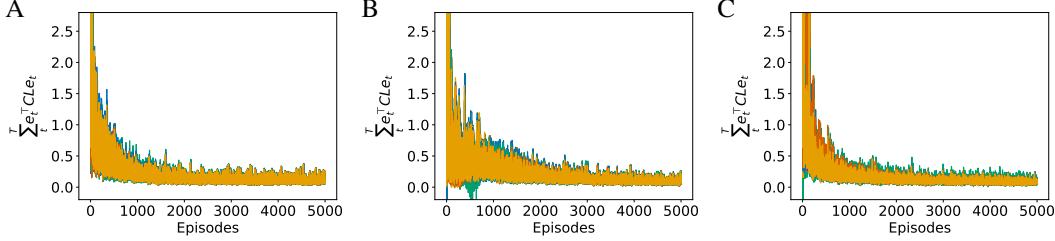
14

Figure S3: **Alignment of weight update $\Delta\theta$ and gradient $g_\theta$.** Average $\sum_t^T e_t^\top CL e_t$ over one episode during system identification as function of episodes for **(A)** 2-d observations (LDS1), **(B)** 3-d observations (LDS2) and **(C)** 3-d observations (LDS2) with an over-representing Bio-OFC that assumes 3 instead of the actual 2 latent dimensions. All 20 runs are shown using different colors.

While we can initialize this way, a further issue is whether the learning rules still minimize the objective at the end of training upon convergence, i.e. near the optimum given by the Kalman filter. Using Eqs. (11-14), we have $g_\theta^\top \Delta\theta = v_{t-\tau} e_t^\top CL e_t v_{t-\tau}^\top$, where $v \in \{\hat{x}, u, e\}$ for $\theta \in \{A, B, L\}$. Hence it suffices if either $CL$ or $LC^{\top+}$ is positive definite. (If the matrices $C$ and $L$ are not square, one of the products won't have full rank, and have eigenvalues that are zero.) For the Kalman filter holds $L = A\Sigma C^\top (C\Sigma C^\top + W)^{-1}$, cf. Eq. (7). It follows that $CL = CA\Sigma C^\top (C\Sigma C^\top + W)^{-1}$. It is reasonable to assume this is positive semi-definite: For now assume $A = I$, then $CA\Sigma C^\top$ is the covariance due to uncertainty of the state, and $(C\Sigma C^\top + W)$ is the total covariance due to uncertainty of the state plus observation noise. The product $CL$ can be considered as a ratio of these covariances, which is close to $I$ for small observation noise. Our learning rules scale the gradients by these covariances. Because $A = I + \mathcal{O}(\Delta t)$ the product $CL$ remains positive definite, as long as the time discretization is not too coarse. Indeed, in the continuous limit of the Kalman-Bucy filter the Kalman gain is simply $L = \Sigma C^\top W^{-1}$.

$g_\theta^\top \Delta\theta > 0$ holds at the end and beginning of training, the latter due to the way we initialize $C$ and $L$. However, what happens throughout learning? When performing stochastic gradient descent only the average update aligns with the negative gradient, whereas individual updates could even increase the objective. Similarly $g_\theta^\top \Delta\theta > 0$ has to hold only on average. Note that $g_\theta^\top \Delta\theta = v_{t-\tau} e_t^\top CL e_t v_{t-\tau}^\top$, where $v \in \{\hat{x}, u, e\}$ for $\theta \in \{A, B, L\}$. We therefore revisited the simulations of Fig. 4 for delay=1 and kept track of $e_t^\top CL e_t$, which needs to be positive on average. While $e_t^\top CL e_t$ was negative for 7.1% of the individual updates for LDS2 (0% for LDS1), the average over one episode $\sum_t^T e_t^\top CL e_t$ was negative for merely 0.036% of the episodes, cf. Fig. S3, and always positive if averaged over multiple episodes. Although we do not present a theoretical derivation to show that $\mathbb{E}[e_t^\top CL e_t] > 0$ the simulations show this is the case.

## E  Consequences of replacing $C^\top$ with L

To obtain local learning rules, Eqs. (13-16), we replaced $C^\top$ in the gradient formulas, Eqs. (11-12), with $L$. Fig. S4 repeats the experiment of Fig. 2 without replacing $C^\top$. The learning rate that minimizes the average MSE over all episodes is smaller, convergence slower, and the average MSE even marginally larger than in Fig. 2. However, the overshooting after the covariance change at 2500 episodes (in Fig. 2A) does not occur. Thus we find that the replacement does not harm performance.

We further noted that replacing $C^\top$ with $L$ corresponds to multiplication of the gradients by $LC^{\top+}$, and we therefore initialize $C$ and $L$ in such a way that $LC^{\top+}$ is positive definite, i.e. its symmetric part has positive real eigenvalues. We investigated the dependence on initial alignment of $C^\top$ and $L$ for the LDS1 experiment and found that the performance does not change in a statistically significant way over a wide range of initial alignments that we considered. In more detail, we learned the Kalman gain $L$ for LDS1 ($C = I$), initializing $L$ as $\begin{pmatrix} 1-a & a \\ a & 1-a \end{pmatrix}$, which has eigenvalues $\lambda_1 = 1 - 2a$ and $\lambda_2 = 1$. If $a = 0$ then $L$ and $C^\top$ are perfectly aligned, for $a = 0.5$ eigenvalue $\lambda_1$ of $L$, and thus of the symmetric part of $LC^{\top+}$, becomes zero. The asymptotic performance (operationally defined as average MSE of the last 100 episodes) for $\lambda_1 > 0.02$ did not differ in a statistical significant way ($p > 0.6$, two sided t-test), but convergence was slower for small but positive $\lambda_1 \gtrsim 0$, cf. Fig. S5. Panel
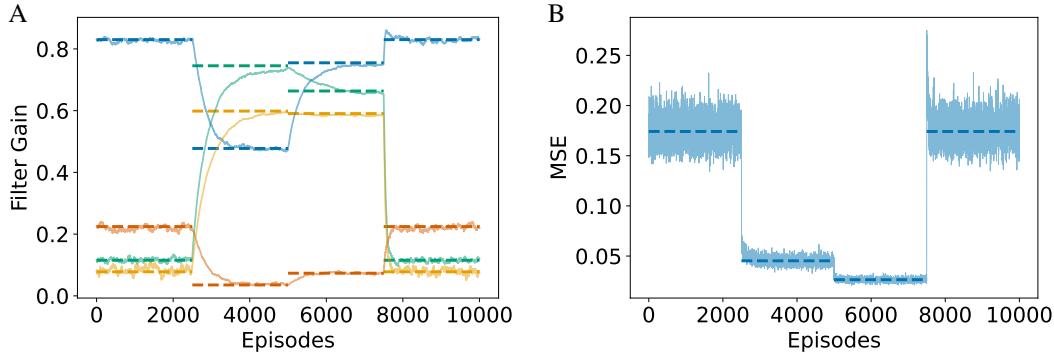
Figure S4: **Adaptive filtering under varying noise levels.** Analogous plots to Fig. 2, but using the non-local SGD learning rule (12) instead of Eq. (14) which replaces $\boldsymbol{C}^\top$ with $\boldsymbol{L}$ to render the learning rule local.



Figure S5: **Dependence on initial alignment of $\boldsymbol{C}^\top$ and $\boldsymbol{L}$.** (A) Average MSE ($\pm$SEM) when the minimal eigenvalue $\lambda_1$ of $\boldsymbol{L}\boldsymbol{C}^{\top+}$ at initialization is varied while the maximal eigenvalue is constant, $\lambda_2 = 1$. Performance does not change in a statistically significant way over a wide range of initial alignments. (B) Convergence to optimal asymptotic performance occurs if $\lambda_1 > 0$ but not if $\lambda_1 \leq 0$.

B shows that for $\lambda_1 \leq 0$ performance did not converge to the optimal asymptotic value. For each value of $\lambda_1$ the learning rate was tuned to minimize the average MSE over all episodes.

Figure S6: **Open-loop training of Bio-OFC on LDS1 for different controller noise levels.** Cost as function of episodes for **(A)** $\sigma = 0.05$, **(B)** $\sigma = 0.1$, **(C)** $\sigma = 0.2$, and **(D)** $\sigma = 0.5$, cf. Fig. 5 for details. Panel C is identical to Fig. 4D.
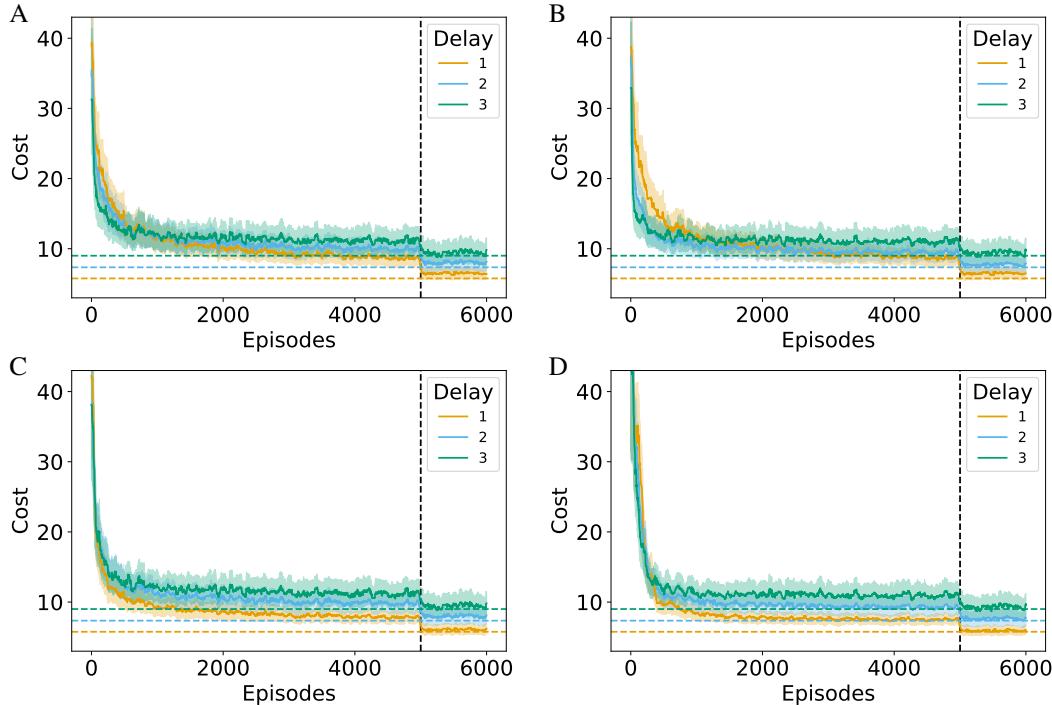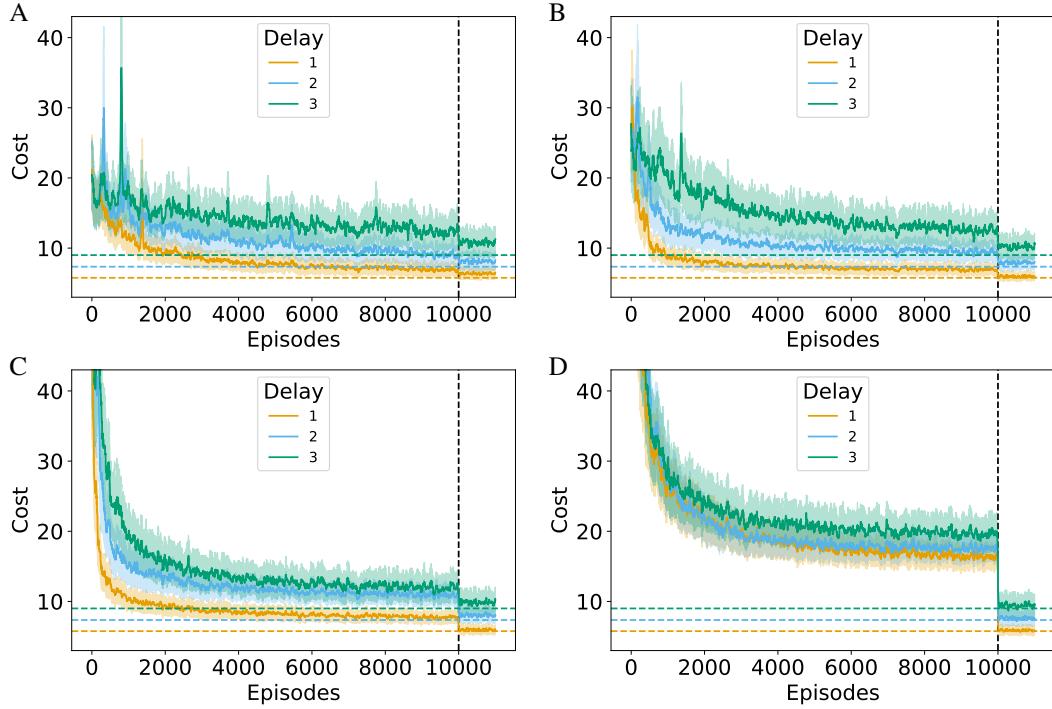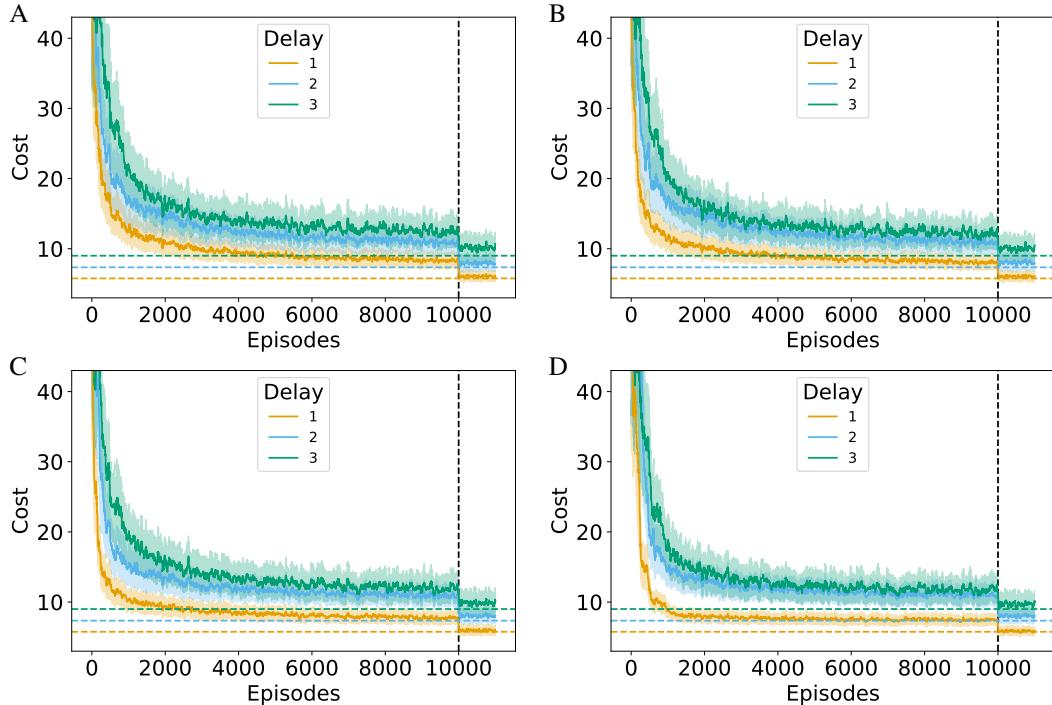


Figure S7: **Open-loop training of Bio-OFC on LDS1 using different momenta in the controller update.** Cost as function of episodes for **(A)** $m = 0$, **(B)** $m = 0.9$, **(C)** $m = 0.99$, and **(D)** $m = 0.9995$, cf. Fig. 4 for details. Panel C is identical to Fig. 4D.

Figure S8: **Closed-loop training of Bio-OFC on LDS1 for different controller noise levels.** Cost as function of episodes for **(A)** $\sigma = 0.05$, **(B)** $\sigma = 0.1$, **(C)** $\sigma = 0.2$, and **(D)** $\sigma = 0.5$, cf. Fig. 5 for details. Panel C is identical to Fig. 5A.



Figure S9: **Closed-loop training of Bio-OFC on LDS1 using different momenta in the controller update.** Cost as function of episodes for **(A)** $m = 0$, **(B)** $m = 0.9$, **(C)** $m = 0.99$, and **(D)** $m = 0.9995$, cf. Fig. 5 for details. Panel C is identical to Fig. 5A.

## F   Learning of a sensory-motor control task

We considered the task of making reaching movements in the presence of externally imposed forces from a mechanical environment [1]. The movements are restricted to a fixed $z$-plane and the 6-dimensional state vector contains the position, velocity and acceleration in $x$ and $y$ direction, $\boldsymbol{x} = (p_x, p_y, v_x, v_y, a_x, a_y)^\top$. In the absence of a force field the system is described by matrices

$$\boldsymbol{A} = \begin{pmatrix} 1 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{pmatrix}, \quad \boldsymbol{B} = \begin{pmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 0 & 0 \\ 1 & 0 \\ 0 & 1 \end{pmatrix}, \quad \boldsymbol{C} = \boldsymbol{I} \tag{S4}$$

For simplicity we follow [2] and assume that the action $\boldsymbol{u}$ is already in Cartesian coordinates as opposed to controlling the torques applied on joints. We leave it for future work to further tighten the connection to biological motor control. We assumed time units of $10\,\mathrm{ms}$, length units of $1\,\mathrm{cm}$ and – in line with experimental data [30] – a measurement delay of $50\,\mathrm{ms}$. Further, the noise covariances and reward matrices were parametrized as

$$\boldsymbol{V} = v\,\mathrm{diag}(1, 1, .1, .1, .01, .01), \quad \boldsymbol{W} = \boldsymbol{V}, \quad \boldsymbol{Q} = \mathrm{diag}(q_1, q_1, q_2, q_2, 0, 0), \quad \boldsymbol{R} = \boldsymbol{I} \tag{S5}$$

where the different scales in $\boldsymbol{V}$ reflect the different numerical scales of position, velocity and acceleration. The forces were computed as a function of the velocity. Application of the force field changes $\boldsymbol{A}$ to

$$\boldsymbol{A} \leftarrow \boldsymbol{A} + f \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & -10.1 & -11.2 & 0 & 0 \\ 0 & 0 & -11.2 & 11.1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \tag{S6}$$

where the numerical values of the non-zero matrix entries were taken from [1]. We chose parameters $v, q_1, q_2, f$ that result in a qualitative match with the experimental data, see Fig. S10 ($v = 10^{-4}$, $q_1 = 10^{-5}$, $q_2 = 0.002$, $f = 0.002$). The optimal filter gain was determined by minimizing the mean squared prediction error over 1000 trajectories. This demonstrated that our network, that approximates OFC, can adequately describe experimental behavior, which is mostly a testimony to the success of OFC. It is standard in linear control theory to put the target state at the origin. This can be achieved by using variables related to the difference between the initial state and the target state. As a result, we can use the same estimator/controller for each reach condition and the different reach conditions correspond to different initial states $x_0$. The more interesting question is whether our plasticity rules Eqs. (13-16) for the filter and Eqs. (21-23) for the controller, capture human performance during the training period. Fig. 6 shows that this is indeed the case and after about 1000 episodes the trajectories are close to straight lines. Fig. S12 repeats this analysis but updates only the system matrices $\boldsymbol{A}, \boldsymbol{B}, \boldsymbol{C}$ and Kalman gain $\boldsymbol{L}$, while keeping the controller weights $\boldsymbol{K}$ fixed. Even switching off learning in the controller yields similar results, thus learning is driven primarily by changes in the estimator. Because the force field alters the system, greater changes in the part that performs system identification, i.e. the estimator, are somewhat to be expected.

To strengthen the connection to biological motor control we reran the reaching task using signal-dependent motor noise in the plant [31] which increases with the magnitude of the control signal. We scaled the amount of noise by the norm of $\boldsymbol{u}$, i.e. we replaced the dynamics $\boldsymbol{x}_{t+1} = \boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t + \boldsymbol{v}_t$ with $\boldsymbol{x}_{t+1} = \boldsymbol{A}\boldsymbol{x}_t + \boldsymbol{B}\boldsymbol{u}_t + |\boldsymbol{u}_t|\boldsymbol{v}_t$. The results, shown in Figs. S13 and S14, are similar to the earlier results obtained with additive noise (Figs. S10 and 6).
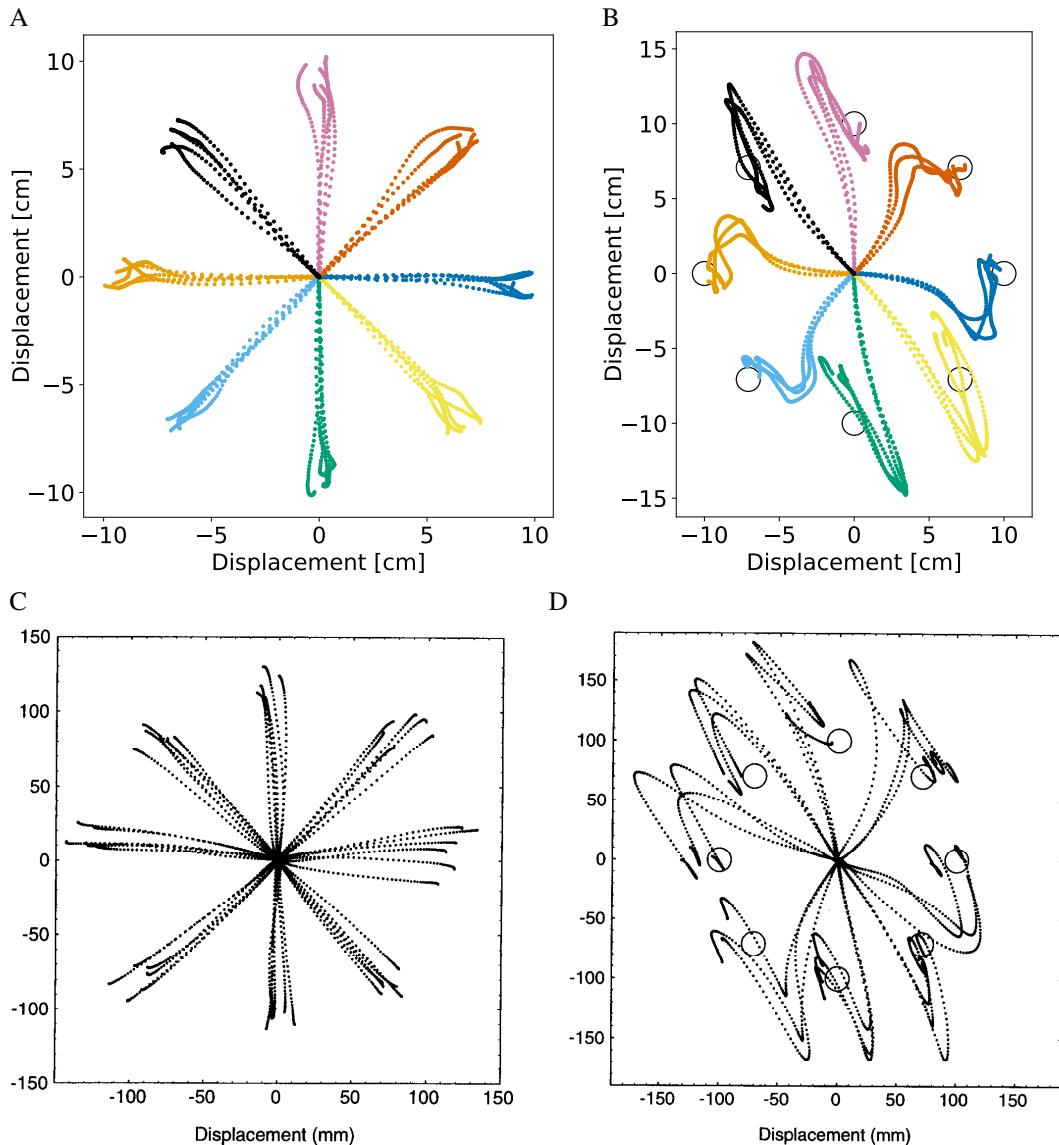
Figure S10: **Optimal feedback control qualitatively captures hand reaching trajectories.** Model trajectories **(A)** in a null force field and **(B)** during initial exposure to a force field. Typical human trajectories **(C)** in a null force field and **(D)** during initial exposure to a force field [1]. Copyright ©1994 Society for Neuroscience.
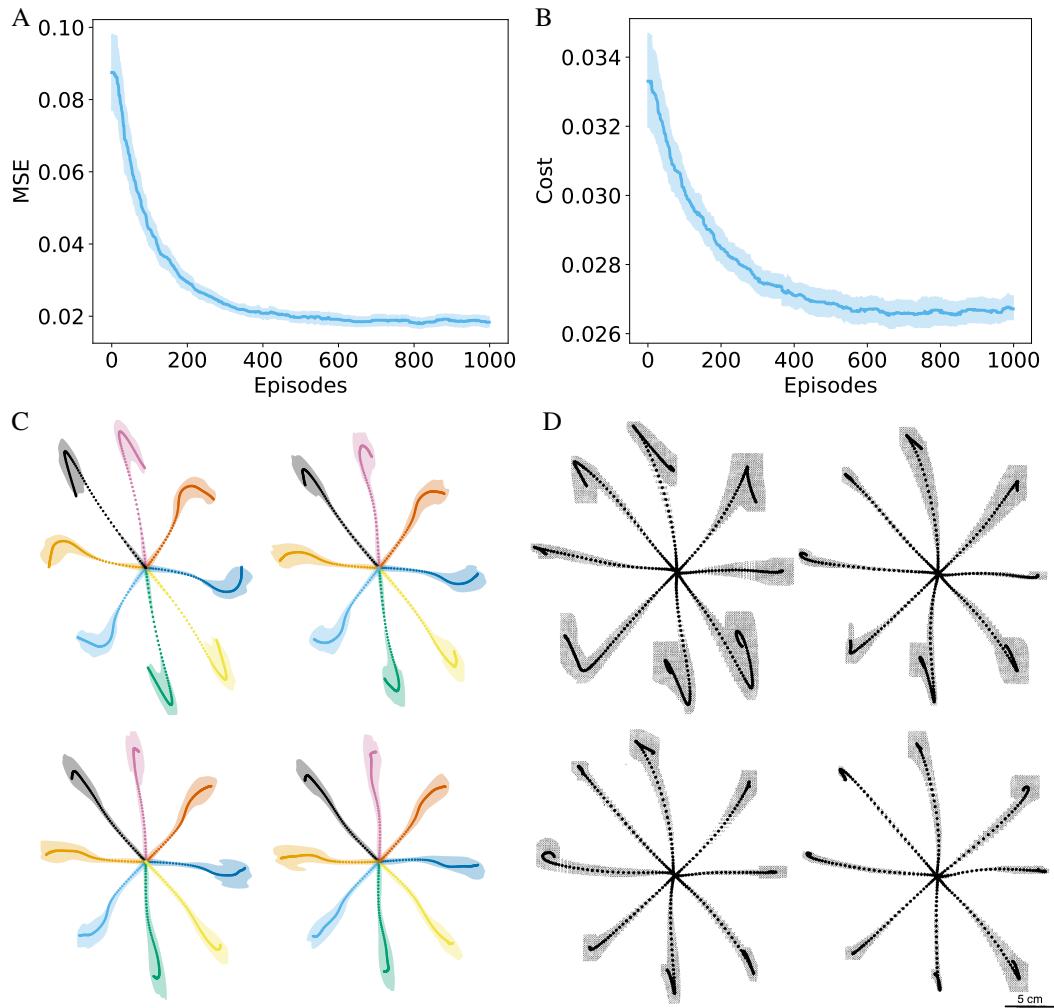
Figure S11: **Learning to adapt to a a force field.** **(A)** Mean squared prediction error and **(B)** cost during training. **(C)** Model trajectories during training. Performance plotted during the first (top left), second (top right), third (bottom left), and final (bottom right) 250 targets. Dots show the mean and are 10 ms apart, shaded area shows a kernel density estimate thresholded at 0.04. **(D)** Averages±SD of human hand trajectories during training [1]. Copyright ©1994 Society for Neuroscience.

Figure S12: **Learning to adapt to a a force field without adaptation of controller weights. (A)** Mean squared prediction error and **(B)** cost during training. **(C)** Model trajectories during training. Performance plotted during the first (top left), second (top right), third (bottom left), and final (bottom right) 250 targets. Dots show the mean and are 10 ms apart, shaded area shows a kernel density estimate thresholded at 0.04. **(D)** Averages±SD of human hand trajectories during training [1]. Copyright ©1994 Society for Neuroscience.
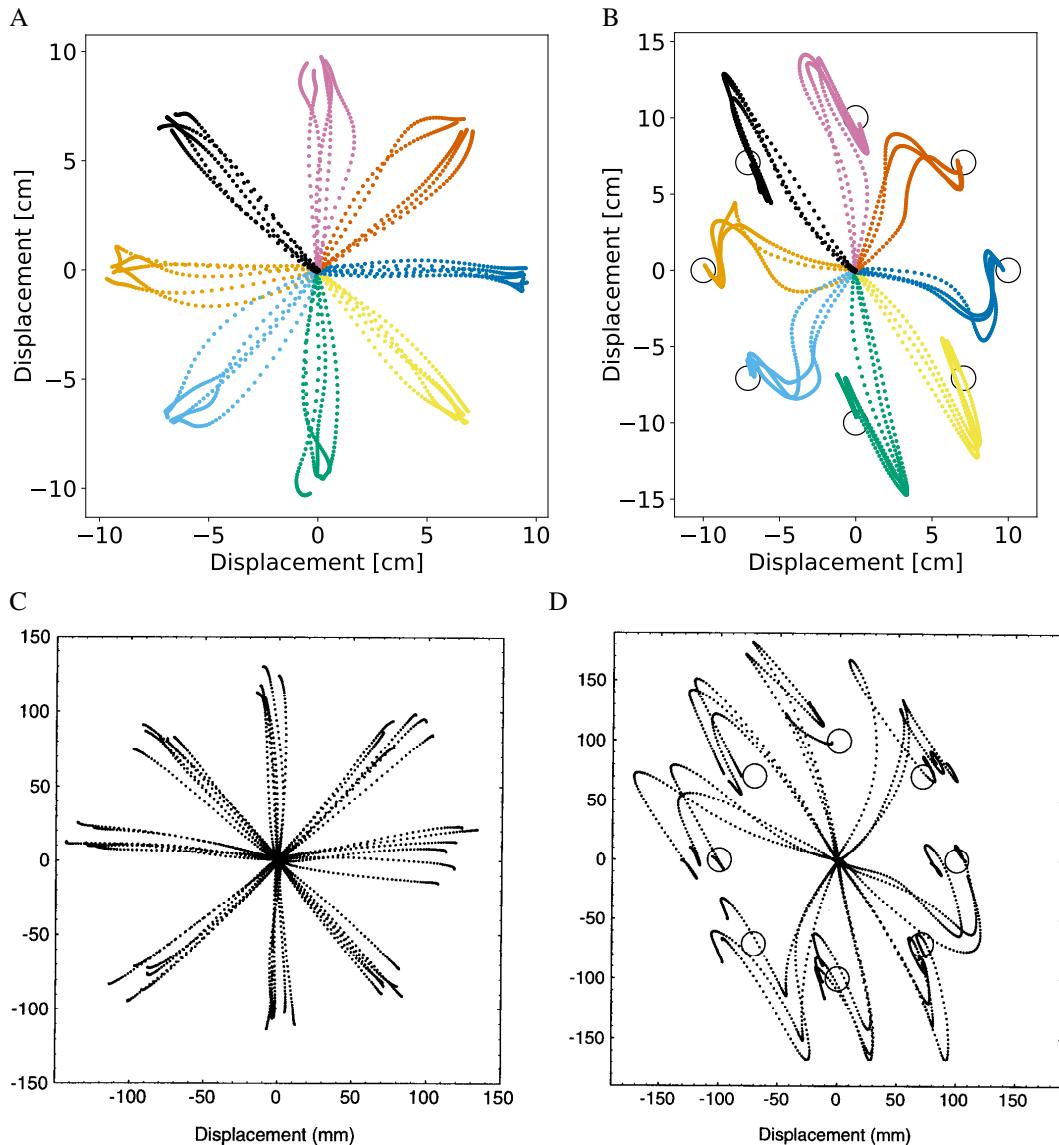
Figure S13: **Optimal feedback control with multiplicative signal-dependent noise qualitatively captures hand reaching trajectories.** Model trajectories **(A)** in a null force field and **(B)** during initial exposure to a force field. Typical human trajectories **(C)** in a null force field and **(D)** during initial exposure to a force field [1]. Copyright ©1994 Society for Neuroscience.
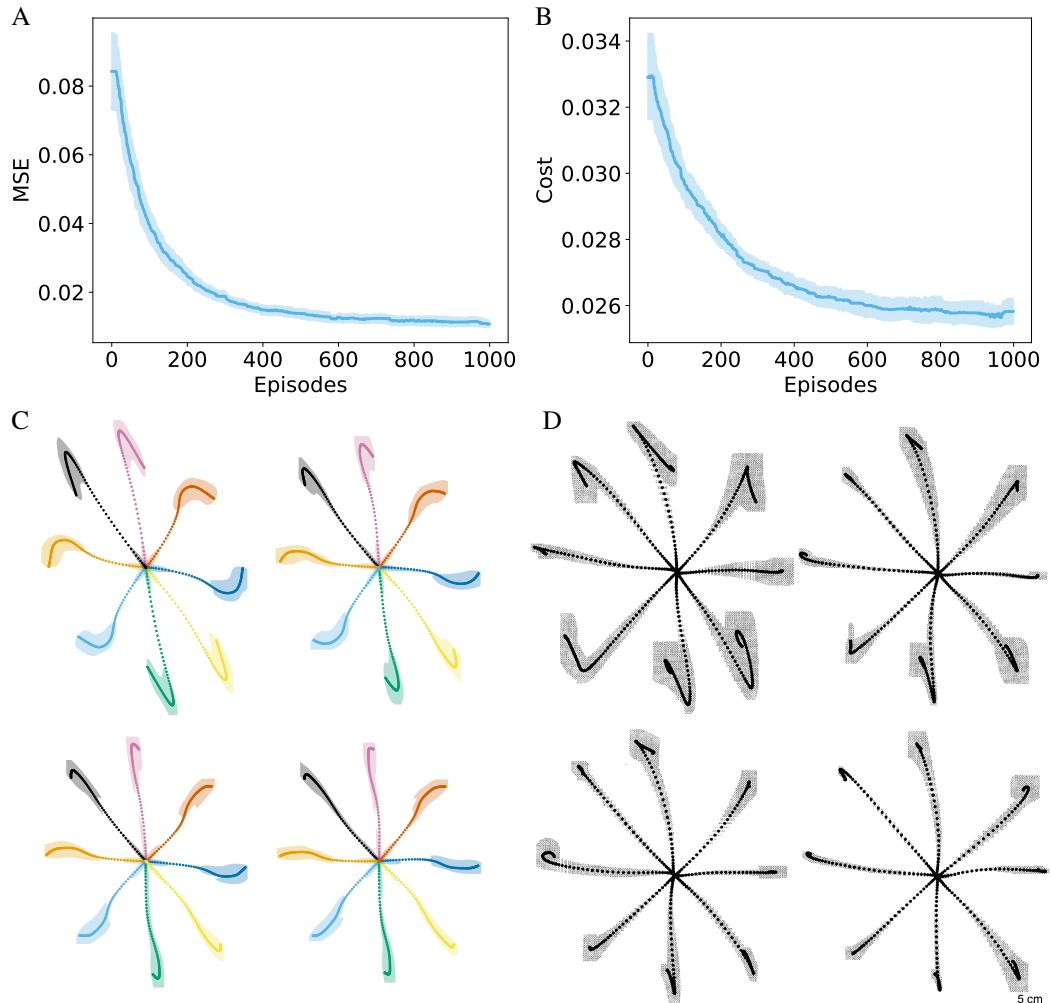
Figure S14: **Learning to adapt to a a force field with multiplicative signal-dependent control noise.** **(A)** Mean squared prediction error and **(B)** cost during training. **(C)** Model trajectories during training. Performance plotted during the first (top left), second (top right), third (bottom left), and final (bottom right) 250 targets. Dots show the mean and are 10 ms apart, shaded area shows a kernel density estimate thresholded at 0.04. **(D)** Averages±SD of human hand trajectories during training [1]. Copyright ©1994 Society for Neuroscience.
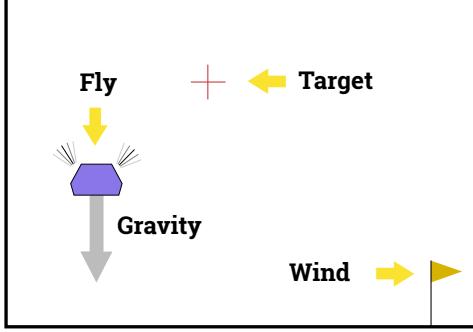
Figure S15: The fly simulation environment.

# G   A simplified fly simulation environment

We designed Gym-Fly, an OpenAI gym [32] environment which simulates flying in a simplified 2-d environment (cf. Fig. S15). We implemented the physical simulations using the Box2D, a 2-d physics engine often used for games [43]. In this environment, the agent controls the flapping frequency of each wing individually. Each wing flap results in an impulse that is up and away from the wing (i.e. direction up and left when flapping the right wing) and for simplicity, we assume that the flapping frequency translates linearly to the magnitude of the impulse imparted at each time-step. Furthermore, we assume that negative flapping frequency results in an impulse in the opposite direction. While this is not realistic in many physical situations, it is necessary for a linear environment which can be described by Eq. (1). The environment has gravity, and wind that randomly increases or decreases at each time-step, up to a maximum value. Furthermore, the system suffers from stochastic noise in that the result of the agent's actions are corrupted before being implemented by the engine. The agent receives sensory stimuli that are composed of its 2-d location and 2-d velocity, as well as the measurement of the wind in the environment. However, these observations are delayed by $100\,\text{ms}$, equivalent to $\tau = 5$ time-steps of the simulation.

The goal of the agent is to fly to a fixed target and stabilize itself against gravity, the environment wind, and stochastic noise in the system. The agent suffers a cost that is proportional to the distance to the target, and the magnitude of the control variables. To verify the flexibility of the Bio-OFC algorithm, we implemented the Gym-Fly environment to deviate from the assumptions used for deriving Bio-OFC in a number of ways. First, the cost suffered by the agent is not quadratic as required by LQR (cf. Eq. (S1)) but is an L1 distance. Second, the system noise that corrupts the control parameters was chosen to be uniform and not Gaussian.

Because of gravity, the description and control of the system with Eqs. (1) and (4), require a bias term. This is because the agent will fall if it stops flapping its wings, that is the point $x = 0$ and $u = 0$ is not a fixed point. Bio-OFC can be easily modified to allow for a bias term. The simplest way to derive the algorithm with a bias is to allow for a component of $\hat{x}$ to be fixed at a constant value, i.e. by replacing $\hat{x} \rightarrow (\hat{x}, 1)$. In a biological setting this bias can be implemented as a thresholding mechanism in the post-synaptic neuron and does not violate the locality and biological plausibility of the learning rules.

We trained Bio-OFC in this environment in a closed-loop setting for a total of 30,000 episodes. The length of the episodes start at 150 time-steps and was linearly increased to 1000 time-steps at episode 500 after which they no longer increase. A video demonstration of how Bio-OFC learns to control this environment is given in `code-gym/gym-fly-demo.mp4`. However, the learning of the bias term, which controls the locations that the agent is stabilized, takes longer to be learned. This is understandable since deviations in the bias term lead to a constant cost at each time-step. However, failure to stabilize the agent against wind or gravity leads to a cost that will diverge in time. We also compared the performance of Bio-OFC to policy gradient. We find that, because of the delay, the agent trained with policy gradient overshoots the target and needs to backtrack, cf. Fig. S16a. However, the agent trained with Bio-OFC flies directly towards the target with no significant overshoot, cf. Fig. S16b.

The code for the Gym-Fly environment (`gym-fly/gym_fly/envs/fly_env.py`) as well as the detailed parameters of the training are provided in the accompanying code (`Bio-OFC gym-fly demo.ipynb`) in the GitHub repository `https://github.com/golkar/bio-ofc-gym`. Installation directions are given in `installation.txt`.



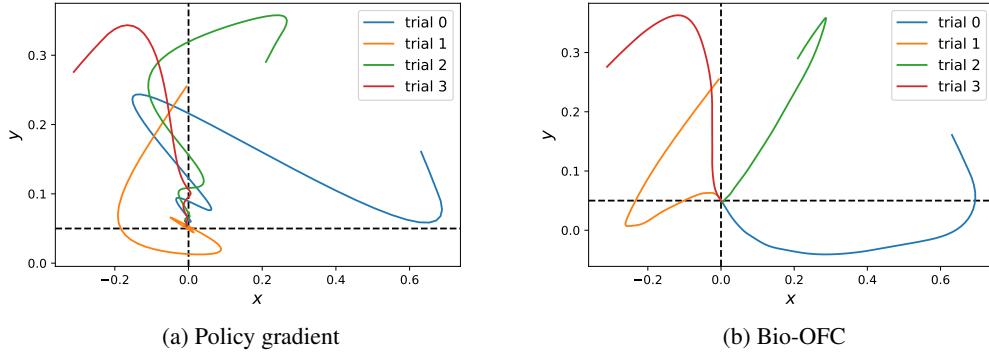(a) Policy gradient

(b) Bio-OFC

Figure S16: Flight trajectory of agents trained with policy gradient (left) and Bio-OFC (right) under the same initial conditions. The agent trained with policy gradient overshoots the target, whereas the agent trained with Bio-OFC has no significant overshoot.