

## Supplementary Information

### 10 Relation between low-pass filter and lookahead

In general, the prospective (or lookahead) voltage  $\check{u}^r$  that enters the activation function “looks into the future” with a time horizon of  $\tau^r$ , that is in general different from  $\tau^m$  (cf. Eqn. 1):

$$\check{u}^r(t) := \left(1 + \tau^r \frac{d}{dt}\right) u(t) . \quad (9)$$

On the other hand, the exponential low-pass filter (LPF) with time constant  $\tau^x$  of some time-dependent quantity  $u(t)$  is given by the average over the past, weighted with an exponential kernel

$$\bar{u}^x(t) := \frac{1}{\tau^x} \int_{-\infty}^t u(t') e^{(t'-t)/\tau^x} dt' . \quad (10)$$

Usually,  $\tau^x$  is either the membrane or the synaptic time constant. To see what happens for a certain neuron in the case of heterogeneous time constants, ( $\tau^m \neq \tau^r$ ), we can compute

$$\frac{d}{dt} \bar{u}^m = \frac{d}{dt} \frac{1}{\tau^m} \int_{-\infty}^t u(t') e^{(t'-t)/\tau^m} dt' \quad (11)$$

$$= \frac{1}{\tau^m} \left( \underbrace{u(t') e^{(t'-t)/\tau^m} \Big|_{t'=t}}_{u(t)} + \int_{-\infty}^t u(t') \frac{\partial}{\partial t} e^{(t'-t)/\tau^m} dt' \right) \quad (12)$$

$$= \frac{1}{\tau^m} \left( u(t) - \underbrace{\frac{1}{\tau^m} \int_{-\infty}^t u(t') e^{(t'-t)/\tau^m} dt'}_{\bar{u}^m(t)} \right) \quad (13)$$

$$= \frac{1}{\tau^m} (u(t) - \bar{u}^m(t)) \quad (14)$$

to obtain the general expression

$$\check{\bar{u}}^m(t) = \left(1 + \tau^r \frac{d}{dt}\right) \bar{u}^m(t) = \bar{u}^m + \frac{\tau^r}{\tau^m} (u(t) - \bar{u}^m(t)) \quad (15)$$

$$= \frac{\tau^r}{\tau^m} u(t) + \frac{\tau^m - \tau^r}{\tau^m} \bar{u}^m(t) . \quad (16)$$

As mentioned in Section 6, different prospective and membrane time constants lead to deviations that persist on the time scale of  $\tau^m$  since the second term in Eqn. 15 is still proportional to the low-pass filtered voltage  $\bar{u}^m$ , which relaxes with the specific time constant  $\tau^m$ . On the other hand, in case of equal time constants  $\tau^m = \tau^r$  it immediately follows that lookahead and LPF are inverse operations:

$$\check{\bar{u}}(t) = \left(1 + \tau \frac{d}{dt}\right) \bar{u}(t) = \bar{u}(t) + \tau \dot{\bar{u}}(t) = u(t) . \quad (17)$$

### 11 Detailed derivation of the neuronal dynamics

Eqn. 3 represents the solution for a stationary energy with respect to the prospective voltage  $\check{u}^m$ . In the following, we show the detailed derivation for an arbitrary component  $i$ :

$$\frac{\partial E}{\partial \check{u}_i^m} = \frac{\partial}{\partial \check{u}_i^m} \frac{1}{2} \sum_{j,k} \|\check{u}_j^m - W_{jk} \varphi(\check{u}_k^m)\|^2 \quad (18)$$

$$= \sum_{j,k,l} [\check{u}_j^m - W_{jk} \varphi(\check{u}_k^m)] \frac{\partial}{\partial \check{u}_i} [\check{u}_j^m - W_{jl} \varphi(\check{u}_l^m)] \quad (19)$$

$$= \sum_{j,k,l} [\check{u}_j^m - W_{jk} \varphi(\check{u}_k^m)] [\delta_{ij} - \delta_{il} W_{jl} \varphi'(\check{u}_l^m)] \quad (20)$$

$$= \check{u}_i^m - \sum_k W_{ik} \varphi(\check{u}_k^m) - \varphi'(\check{u}_i^m) \sum_{j,k} W_{ij}^T [\check{u}_j^m - W_{jk} \varphi(\check{u}_k^m)] \quad (21)$$

and therefore

$$0 = \frac{\partial E}{\partial \check{u}_i^m} \implies \tau^m \dot{u}_i = -u_i + \sum_k W_{ik} \varphi(\check{u}_k^m) + \varphi'(\check{u}_i^m) \sum_{j,k} W_{ij}^T [\check{u}_j^m - W_{jk} \varphi(\check{u}_k^m)] \quad (22)$$

### 11.1 Neuronal dynamics with synaptic filtering

To include synaptic filtering in our theory, we introduce an additional LPF as in Eqn. 10 with time constant  $\tau^s$ . For this, it is sufficient to replace firing rates with filtered rates,  $\varphi \rightarrow \bar{\varphi}^s$ , in the total energy  $E$ :

$$E(\check{\mathbf{u}}^m) := \frac{1}{2} \|\check{\mathbf{u}}^m - \mathbf{W} \bar{\boldsymbol{\varphi}}^s(\check{\mathbf{u}}^m) - \mathbf{b}\|^2 + \beta \mathcal{L}(\check{\mathbf{u}}^m). \quad (23)$$

Deriving the neuronal dynamics from a vanishing gradient  $\nabla_{\check{\mathbf{u}}^m} E$  as before now yields

$$\tau^m \dot{\mathbf{u}} = -\mathbf{u} + \mathbf{W} \bar{\boldsymbol{\varphi}}^s(\check{\mathbf{u}}^m) + \mathbf{b} + \mathbf{e}, \quad (24)$$

where the error term now reads

$$\mathbf{e} = \bar{\boldsymbol{\varphi}}^s(\check{\mathbf{u}}^m) \mathbf{W}^T [\check{\mathbf{u}}^m - \mathbf{W} \bar{\boldsymbol{\varphi}}^s(\check{\mathbf{u}}^m) - \mathbf{b}]. \quad (25)$$

These are essentially the same equations as before, just with the rates replaced with their filtered version.

## 12 General formulation for arbitrary connectivity functions

Here we consider a generalization of the energy function from the main manuscript that includes arbitrary ‘‘connectivity functions’’  $f$  with parameters  $\boldsymbol{\theta}$ :

$$E(\check{\mathbf{u}}^m) := \frac{1}{2} \|\check{\mathbf{u}}^m - \mathbf{f}(\boldsymbol{\varphi}(\check{\mathbf{u}}^m), \boldsymbol{\theta})\|^2 + \beta \mathcal{L}(\check{\mathbf{u}}^m). \quad (26)$$

Again, we derive neuron dynamics by requiring  $\nabla_{\check{\mathbf{u}}^m} E(\check{\mathbf{u}}^m) = 0$ . For simplicity, we compute it element wise, shown here for a neuron that does not directly contribute to the loss  $\mathcal{L}$ :

$$\begin{aligned} \frac{\partial E(\check{\mathbf{u}}^m)}{\partial \check{u}_i^m} &= \check{u}_i^m - f_i(\boldsymbol{\varphi}(\check{\mathbf{u}}^m), \boldsymbol{\theta}) + \sum_j \frac{\partial}{\partial \check{u}_i^m} \frac{1}{2} (\check{u}_j^m - f_j(\boldsymbol{\varphi}(\check{\mathbf{u}}^m), \boldsymbol{\theta}))^2 \\ &= \check{u}_i^m - f_i(\boldsymbol{\varphi}(\check{\mathbf{u}}^m), \boldsymbol{\theta}) - \sum_j \frac{\partial \varphi_i}{\partial \check{u}_i^m} \frac{\partial f_j(\boldsymbol{\varphi}(\check{\mathbf{u}}^m), \boldsymbol{\theta})}{\partial \varphi_i} (\check{u}_j^m - f_j(\boldsymbol{\varphi}(\check{\mathbf{u}}^m), \boldsymbol{\theta})) \end{aligned} \quad (27)$$

From this we can, for example, obtain the neuron dynamics described in the main manuscript by choosing a linear connectivity function  $f_i = \sum_j w_{ij} \varphi(\check{u}_j^m) + b_i$ .

## 13 Simulation details

### 13.1 Gradient-based model

For Fig. 1c we consider the neuron dynamics from the main manuscript, but replace prospective membrane potentials with their instantaneous version. Furthermore, we consider a squared loss with some fixed target  $t^*$ . To avoid artificially introducing an additional mismatch problem, we add an exponential low-pass filter to the error terms; this prevents the model to reduce weights to zero in the absence of a teacher due to a continuous mismatch between instantaneous bottom-up predictions and slow neuronal responses. This results in the following dynamics for the two neurons:

$$\tau^m \dot{u}_1 = -u_1 + w_1 r_{\text{in}} + \varphi'(u_1) w_2 (u_2 - w_2 \bar{\varphi}(u_1)) \quad (28)$$

$$\tau^m \dot{u}_2 = -u_2 + w_2 \varphi(u_1) + \beta (t^* - u_2) \quad (29)$$

As in the main manuscript, plasticity is defined as stochastic gradient descent on the energy; as above, we consider low-pass filtered variants of inputs:

$$\Delta w_i = \eta (u_i - W \bar{r}_{i-1}) \bar{r}_{i-1} \quad (30)$$

Here we use linear activation functions and the following parameters:  $dt = 0.001$ ,  $\tau^m = 10\text{ms}$ ,  $\beta \in \{0, 0.9\}$ ,  $\eta_W = 0.0005$ .

### 13.2 Numerical implementation

In order to carry out our simulations we had to discretize the differential equations, which we state here again for clarity:

$$\tau^m \dot{\mathbf{u}}_\ell(t) = -\mathbf{u}_\ell(t) + \mathbf{W}_\ell \varphi(\check{\mathbf{u}}_\ell^m(t)) + \varphi'(\check{\mathbf{u}}_\ell^m(t)) \mathbf{W}_{\ell+1}^T [\check{\mathbf{u}}_{\ell+1}^m(t) - \mathbf{W}_{\ell+1} \varphi(\check{\mathbf{u}}_\ell^m(t))] \quad (31)$$

$$\dot{\mathbf{W}}_\ell(t) = \eta [\check{\mathbf{u}}_\ell^m(t) - \mathbf{W}_\ell \varphi(\check{\mathbf{u}}_{\ell-1}^m(t))] \varphi(\check{\mathbf{u}}_{\ell-1}^m(t)) . \quad (32)$$

First, observe that  $\dot{\mathbf{u}}$  appears on both sides of the first equation: explicitly on the left hand side, and implicitly on the right, in the definition of  $\check{\mathbf{u}}^m(t) = \mathbf{u}(t) + \tau^m \dot{\mathbf{u}}(t)$ . To resolve this circular dependency, we define  $\check{\mathbf{u}}^m(t + dt) = \mathbf{u}(t) + \tau^m \dot{\mathbf{u}}(t)$ , which works well for small enough  $dt$ .

We first consider the neuronal update and use forward Euler to rewrite the derivative as a finite difference with time step  $dt$

$$\tau^m \frac{\mathbf{u}_\ell(t + dt) - \mathbf{u}_\ell(t)}{dt} = -\mathbf{u}_\ell(t) + \mathbf{W}_\ell \varphi(\check{\mathbf{u}}_\ell^m(t)) + \mathbf{e}_\ell(t) \quad (33)$$

with

$$\mathbf{e}_\ell(t) := \varphi'(\check{\mathbf{u}}_\ell^m(t)) \mathbf{W}_{\ell+1}^T [\check{\mathbf{u}}_{\ell+1}^m(t) - \mathbf{W}_{\ell+1} \varphi(\check{\mathbf{u}}_\ell^m(t))] \quad (34)$$

and solve for  $\mathbf{u}_\ell(t + dt)$  to obtain

$$\mathbf{u}_\ell(t + dt) = \mathbf{u}_\ell(t) + dt \Delta \mathbf{u}_\ell(t) , \quad (35)$$

where

$$\Delta \mathbf{u}_\ell(t) = \frac{1}{\tau^m} [-\mathbf{u}_\ell(t) + \mathbf{W}_\ell \varphi(\check{\mathbf{u}}_\ell^m(t)) + \mathbf{e}_\ell(t)] . \quad (36)$$

Similarly, we use forward Euler for weight dynamics:

$$\frac{\mathbf{W}_\ell(t + dt) - \mathbf{W}_\ell(t)}{dt} = \eta [\check{\mathbf{u}}_\ell^m(t) - \mathbf{W}_\ell \varphi(\check{\mathbf{u}}_{\ell-1}^m(t))] \varphi(\check{\mathbf{u}}_{\ell-1}^m(t)) \quad (37)$$

and

$$\mathbf{W}_\ell(t + dt) = \mathbf{W}_\ell(t) + dt \Delta \mathbf{W}_\ell(t) , \quad (38)$$

with

$$\Delta \mathbf{W}_\ell(t) = \eta [\check{\mathbf{u}}_\ell^m(t) - \mathbf{W}_\ell \varphi(\check{\mathbf{u}}_{\ell-1}^m(t))] \varphi(\check{\mathbf{u}}_{\ell-1}^m(t)) \quad (39)$$

For both  $\mathbf{e}_\ell(t)$  and  $\Delta \mathbf{W}_\ell(t)$  it is crucial to combine the  $\check{\mathbf{u}}^m(t)$  with the input  $\mathbf{W}_\ell \varphi(\check{\mathbf{u}}_{\ell-1}^m(t))$  that was also used in computing  $\check{\mathbf{u}}^m(t)$ . Pseudo-code for our vanilla implementation can be found in Algorithm 1.

---

#### Algorithm 1 Pseudo-code for the multi-layer implementation of Latent Equilibrium (LE)

---

```

1: for all layers  $\ell$  from 1 (input) to  $N$  (output) do
2:    $\mathbf{e}_\ell(t) \leftarrow (1 - \delta_{\ell N}) \varphi'(\check{\mathbf{u}}_\ell^m(t)) \mathbf{W}_{\ell+1}^T(t) [\check{\mathbf{u}}_{\ell+1}^m(t) - \mathbf{W}_{\ell+1}(t) \varphi(\check{\mathbf{u}}_\ell^m(t))] + \delta_{\ell N} \mathbf{e}^{\text{trg}}(t)$ 
3:    $\Delta \mathbf{u}_\ell(t) \leftarrow \tau^{-1} [-\mathbf{u}_\ell(t) + \mathbf{W}_\ell(t) \mathbf{r}_{\ell-1}(t) + \mathbf{e}_\ell(t)]$ 
4:    $\mathbf{u}_\ell(t + dt) \leftarrow \mathbf{u}_\ell(t) + dt \Delta \mathbf{u}_\ell(t)$ 
5:    $\check{\mathbf{u}}_\ell(t + dt) \leftarrow \mathbf{u}_\ell(t) + \tau \Delta \mathbf{u}_\ell(t)$ 
6:   if synaptic plasticity then
7:      $\Delta \mathbf{W}_\ell(t) \leftarrow \eta \mathbf{e}_\ell(t) \cdot \varphi(\check{\mathbf{u}}_{\ell-1}^m(t))^T$ 
8:      $\mathbf{W}_\ell(t + dt) \leftarrow \mathbf{W}_\ell(t) + dt \Delta \mathbf{W}_\ell(t)$ 
9:   end if
10: end for

```

---

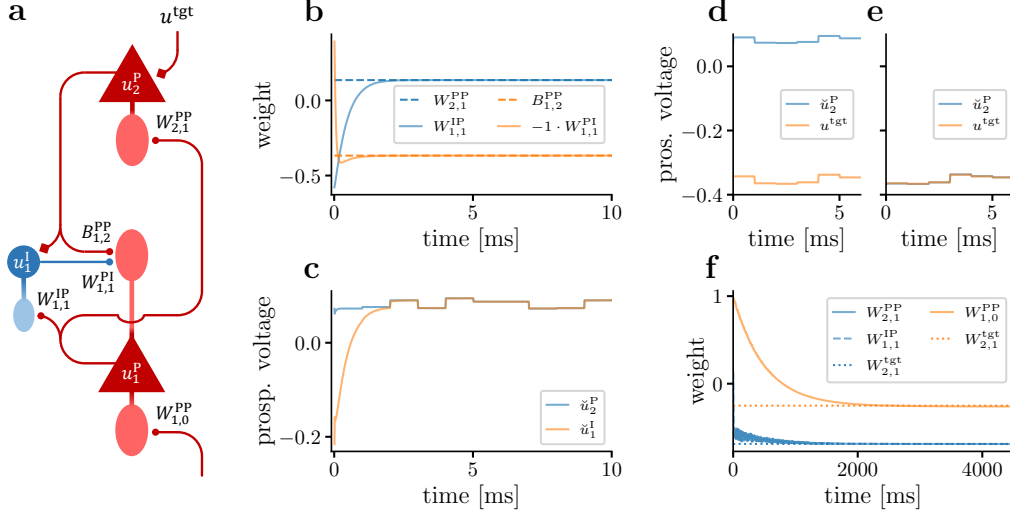
## 14 Microcircuit details

The somatic membrane potential of hidden layer pyramidal cells, interneurons, and top-layer pyramidal cells is described by the following differential equations:

$$C_m \dot{\mathbf{u}}_\ell^P = g_l (E_l - \mathbf{u}_\ell^P) + g^{\text{bas}} (\mathbf{v}_\ell^{\text{bas}} - \mathbf{u}_\ell^P) + g^{\text{api}} (\mathbf{v}_\ell^{\text{api}} - \mathbf{u}_\ell^P) , \quad (40)$$

$$C_m \dot{\mathbf{u}}_\ell^I = g_l (E_l - \mathbf{u}_\ell^I) + g^{\text{den}} (\mathbf{v}_\ell^{\text{den}} - \mathbf{u}_\ell^I) + i^{\text{nudge, I}} , \quad (41)$$

$$C_m \dot{\mathbf{u}}_N^P = g_l (E_l - \mathbf{u}_N^P) + g^{\text{bas}} (\mathbf{v}_N^{\text{bas}} - \mathbf{u}_N^P) + i^{\text{nudge, tgt}} . \quad (42)$$



**Figure 5: Learning to mimic a teacher microcircuit with LE.** (a) Microcircuit architecture following [1]. (b) Learning of the lateral weights  $W_{1,1}^{IP}$  and  $W_{1,1}^{PI}$  to implement the self-predicting state. (c) Prospective membrane voltages during learning of the self-predicting state where (in absence of a target) the top-down activity is matched by the activity of the interneuron. (d, e) Comparison between the prospective membrane voltage  $\tilde{u}_2^P$  of the output pyramidal neuron and the target voltage  $u^{\text{tgt}}$  before (d) and after (e) training. (f) Weight evolution during learning.

Target signals to interneurons and top-layer pyramidal neurons are modeled as a conductance-based input to the respective somatic compartments:

$$i^{\text{nudge, I}} = g^{\text{nudge, I}} (\mathbf{u}_{\ell+1}^P - \mathbf{u}_\ell^I) , \quad (43)$$

$$i^{\text{nudge, tgt}} = g^{\text{nudge, tgt}} (\mathbf{u}^{\text{tgt}} - \mathbf{u}_N^P) . \quad (44)$$

The target signal for the top-layer pyramidal is determined by the training set. For the interneurons, the somatic membrane potentials of the pyramidal neurons in the layer above serve as targets. The membrane potentials of the dendritic compartments instantaneously follow their inputs:

$$\mathbf{v}_\ell^{\text{bas}} = \mathbf{W}_{\ell, \ell-1}^{\text{PP}} \varphi(\mathbf{u}_{\ell-1}^P) , \quad (45)$$

$$\mathbf{v}_\ell^{\text{api}} = \mathbf{B}_{\ell, \ell+1}^{\text{PP}} \varphi(\mathbf{u}_{\ell+1}^P) + \mathbf{W}_{\ell, \ell}^{\text{PI}} \varphi(\mathbf{u}_\ell^I) , \quad (46)$$

$$\mathbf{v}_\ell^{\text{den}} = \mathbf{W}_{\ell, \ell}^{\text{IP}} \varphi(\mathbf{u}_\ell^P) . \quad (47)$$

All synapses except the top-down connections are plastic. The learning rules are described by

$$\dot{\mathbf{W}}_{\ell, \ell-1}^{\text{PP}} = \eta_\ell^{\text{PP}} \left[ \varphi(\mathbf{u}_\ell^P) - \varphi\left(\frac{g^{\text{bas}}}{g_1 + g^{\text{bas}} + g^{\text{api}}} \mathbf{v}_\ell^{\text{bas}}\right) \right] \varphi(\mathbf{u}_{\ell-1}^P) , \quad (48)$$

$$\dot{\mathbf{W}}_{\ell, \ell}^{\text{IP}} = \eta_\ell^{\text{IP}} \left[ \varphi(\mathbf{u}_\ell^I) - \varphi\left(\frac{g^{\text{den}}}{g_1 + g^{\text{den}}} \mathbf{v}_\ell^{\text{den}}\right) \right] \varphi(\mathbf{u}_\ell^P) , \quad (49)$$

$$\dot{\mathbf{W}}_{\ell, \ell}^{\text{PI}} = \eta_\ell^{\text{PI}} \left[ -\mathbf{v}_\ell^{\text{api}} \right] \varphi(\mathbf{u}_\ell^I) . \quad (50)$$

Here, the weights of the top-down connections  $B^{\text{PP}}$  are static and random.

The differential equations for the somatic membrane potentials of all neuron types can be rewritten in a simpler form which also increases their numerical stability:

$$C_m \dot{\mathbf{u}} = \frac{1}{\tau_{\text{eff}}} (\mathbf{u}^{\text{eff}} - \mathbf{u}) , \quad (51)$$

$$\tau_{\text{eff}} = \frac{C_m}{g_1 + g^{\text{bas/den}} + g^{\text{api/nudge}}} . \quad (52)$$

Here,  $\mathbf{u}^{\text{eff}}$  is the effective reversal potential defined as:

$$\mathbf{u}_\ell^{\text{eff,P}} = \frac{g_1 E_1 + g^{\text{bas}} \mathbf{v}_\ell^{\text{bas}} + g^{\text{api}} \mathbf{v}_\ell^{\text{api}}}{g_1 + g^{\text{bas}} + g^{\text{api}}}, \quad (53)$$

$$\mathbf{u}_\ell^{\text{eff,I}} = \frac{g_1 E_1 + g^{\text{den}} \mathbf{v}_\ell^{\text{den}} + g^{\text{nudge,I}} \mathbf{u}_{\ell+1}^{\text{P}}}{g_1 + g^{\text{den}} + g^{\text{nudge,I}}}, \quad (54)$$

$$\mathbf{u}_N^{\text{eff,P}} = \frac{g_1 E_1 + g^{\text{bas}} \mathbf{v}_N^{\text{bas}} + g^{\text{nudge,tgt}} \mathbf{u}^{\text{tgt}}}{g_1 + g^{\text{bas}} + g^{\text{nudge,tgt}}} \quad (55)$$

if a target is provided. If no target is provided to the top-layer pyramidal neurons we assume  $\mathbf{u}^{\text{tgt}} = \mathbf{u}_N^{\text{eff,P}}$  and the above equation simplifies to

$$\mathbf{u}_N^{\text{eff,P}} = \frac{g_1 E_1 + g^{\text{bas}} \mathbf{v}_N^{\text{bas}}}{g_1 + g^{\text{bas}}}. \quad (56)$$

We include LE in the dendritic microcircuit by two simple modifications. First, the output rate of the neurons must depend on the prospective voltage:  $\varphi(\mathbf{u}) \rightarrow \varphi(\check{\mathbf{u}})$ . Note that this includes also the rates in the calculation of dendritic membrane potentials (Eqns. 45 to 47) as well as the plasticity rules (Eqns. 48 to 50). Secondly, the nudging for the interneurons must depend on the prospective voltage of the pyramidal neurons above:

$$\mathbf{u}_\ell^{\text{eff,I}} = \frac{g_1 E_1 + g^{\text{den}} \mathbf{v}_\ell^{\text{den}} + g^{\text{nudge,I}} \check{\mathbf{u}}_{\ell+1}^{\text{P}}}{g_1 + g^{\text{den}} + g^{\text{nudge,I}}}. \quad (57)$$

For the simulation of the cortical networks shown in Fig. 3 and Fig. 5, we use the Euler integration method. Similarly to the LE networks without microcircuit connectivity, we break the circular dependency of  $\check{\mathbf{u}}^{\text{m}}$  in an Euler integration step by defining  $\check{\mathbf{u}}^{\text{m}}$  as a function of previous time steps (see Section 13.2).

Learning is split into two stages: first, the learning of the so-called self-predicting state and afterwards the learning of the actual task. The self-predicting state describes a configuration of weights in which, in the absence of target signals provided to the last layer, apical dendrites are always at rest and the somatic membrane potentials of the interneurons match the membrane potentials of the pyramidal neurons in the layer above. In this state, the network is able to correctly transport errors induced by the target signal to the apical compartments of the lower layer neurons.

Here we demonstrate, for a single microcircuit, the learning of the self-predicting state from a random initialization of weights, by presenting the network with random inputs, no targets to the output layer and using the learning rules given above with  $\eta^{\text{PP}} = 0$  and  $\eta^{\text{PI/IP}} \neq 0$  (Fig. 5 b, c). After the self-predicting state is learned, the network is taught to reproduce the input-output relationship produced by a teacher network (Fig. 5 d-f). For the learning of the task we set the learning rates to  $\eta^{\text{PI}} = 0$  and  $\eta^{\text{PP/IP}} \neq 0$ . In the main manuscript (Fig. 3) we initialize the weights with

$$\mathbf{W}_{1,1}^{\text{IP}} = \frac{g^{\text{bas}} (g_1 + g^{\text{den}})}{g^{\text{den}} (g_1 + g^{\text{bas}})} \mathbf{W}_{2,1}^{\text{PP}} \quad \text{and} \quad (58)$$

$$\mathbf{W}_{1,1}^{\text{PI}} = -\mathbf{B}_{1,2}^{\text{PP}}, \quad (59)$$

thereby skipping the first learning stage and initializing the network directly in the self-predicting state. The full set of parameters used in Fig. 3 and Fig. 5 can be found in Section 15.2.

## 15 Parameters

### 15.1 Parameters used for classification experiments shown in Fig. 2

Table 1 lists all the parameters we used for the experiments shown in Fig. 2. This includes HIGGS and MNIST experiments with fully connected (FC) architectures in Fig. 2b and c as well as MNIST and CIFAR-10 experiments employing convolutional networks (ConvNets).

Standard artificial neural networks (ANNs) were trained with classical backpropagation (BP) using the same network topologies but with cross-entropy (CE) loss instead of the mean squared error (MSE) loss used for the LE experiments.

**Table 1:** Neuron, network and training parameters used to produce the results shown in Fig. 2.

| Symbol                           | Parameter name                        | Fig. 2a              | Fig. 2b            | Fig. 2c                                    | Fig. 2d                          | Fig. 2e                          |
|----------------------------------|---------------------------------------|----------------------|--------------------|--|----------------------------------|----------------------------------|
| <b>Neuron parameters</b>         |                                       |                      |                    |  |                                  |                                  |
| $\tau^m$ [ms]                    | membrane time constant                | 10                   | 20                 | 10   | 10                               | 10                               |
| $\tau^r$ [ms]                    | prospective time constant             | 10                   | 20                 | 10   | 10                               | 10                               |
| $\tau^s$ [ms]                    | synaptic time constant                | 0                    | 0                  | 0  | 0                                | 0                                |
| $\varphi_\ell$                   | activation                            | tanh                 |                    | hard sigmoid <sup>1</sup>                  |                                  |                                  |
| $\varphi_N$                      | output activation                     | linear               |                    |  |                                  |                                  |
| <b>Network parameters</b>        |                                       |                      |                    |  |                                  |                                  |
|                                  | architecture                          | FC                   | FC                 | FC   | LeNet-5                          | LeNet-5                          |
|                                  | input size                            | 50                   | 784                | 28   | $28 \times 28 \times 1$          | $32 \times 32 \times 3$          |
|                                  | hidden layer size                     | 30                   | 300                | 300  | $C(5 \times 5) \times 20 - MP^2$ | $C(5 \times 5) \times 50 - MP^2$ |
|                                  |                                       |                      | 100                | 300  | $C(5 \times 5) \times 50 - MP^2$ | 500                              |
|                                  | output layer size                     | 1                    | 10                 | 1  | 10                               |                                  |
| $\beta$                          | nudging strength                      | 0.1                  |                    |  |                                  |                                  |
| $\mathcal{L}$                    | loss                                  | MSE                  |                    |  |                                  |                                  |
|                                  | initial weights & biases              | uniform <sup>2</sup> |                    | $\sim \mathcal{N}(\mu = 0, \sigma = 0.05)$ |                                  |                                  |
| $\eta_{w,b}$ [ms <sup>-1</sup> ] | learning rate                         | 0.25                 | $128 \times 0.125$ | $64 \times 0.125$                          | $128 \times 0.125$               | $64 \times 0.125$                |
|                                  | layerwise $\eta$ factors <sup>3</sup> | –                    | 1, .2, .1          | 1, .2, .1, .1                              | 1, .2, .1                        | 1, .2, .2, .2, .1                |
| <b>Training parameters</b>       |                                       |                      |                    |  |                                  |                                  |
| $dt$ [ms]                        | temporal resolution                   | 0.001                | 0.01               | 0.1  | 0.1                              | 0.1                              |
| $T_{\text{pres}}$                | presentation time                     | 1 dt                 | 100 dt             | 20 dt                                      | 100 dt                           | 50 dt                            |
|                                  |                                       | = 0.001 ms           | = 1 ms             | = 2 ms                                     | = 10 ms                          | = 5 ms                           |
|                                  | batch size                            | 1                    | 512                | 128  | 512                              | 128                              |
|                                  | # training epochs                     | –                    |                    |  | 100                              |                                  |
|                                  | # train samples                       | –                    | 50000              | 40000                                      | 50000                            | 342000                           |
|                                  | # validation samples                  | –                    | 10000              | 10000                                      | 10000                            | 18000                            |
|                                  | # test samples                        | –                    | 10000              | 10000                                      | 10000                            | 40000                            |
|                                  | # seeds                               | –                    | 10                 | 9  | 9                                | 9                                |

<sup>1</sup> By “hard sigmoid” we mean the piecewise linear function  $\varphi(x)$  that is obtained clipping a rectified linear unit (ReLU) to  $[0, 1]$

$$\varphi(x) = x\theta(x) - (x-1)\theta(x-1) = \begin{cases} 0 & \text{if } x \leq 0 \\ 1 & \text{if } x \geq 1 \\ x & \text{else} \end{cases}, \quad \varphi'(x) = \begin{cases} 1 & \text{if } x \in [0, 1] \\ 0 & \text{else} \end{cases}$$

where  $\theta(x)$  denotes the Heaviside step function.

<sup>2</sup> PyTorch defaults

<sup>3</sup> These factors scale the learning rate  $\eta$  for each layer independently.

<sup>4</sup> C and MP indicate convolutional and max pooling layers, respectively.

Also, we used a ReLU function for the hidden layer activation of the ANNs instead of the hard sigmoid activation that was used for the LE simulations. Furthermore, the BP results for the HIGGS dataset were produced using different activation functions for both hidden and output layers, namely tanh and sigmoidal, respectively.

To perform the simulations without prospective coding shown in Fig. 2b we set  $\tau^r = 0$  ms,  $\tau^m = 10$  ms and used a temporal resolution of  $dt = 1$  ms to obtain reasonable presentation times of multiple  $\tau^m$  that allow for relaxation. That is to say  $T_{\text{pres}} = 200 dt = 200$  ms  $= 20 \tau^m$  for the purple curve in Fig. 2 compared to presentation times of  $T_{\text{pres}} = 1$  ms  $= 0.05 \tau^m$  used for the LE simulations that employ the prospective coding allowing for much smaller presentation times.

## 15.2 Parameters used for the experiments shown in Fig. 3 and Fig. 5

**Table 2:** Neuron, network and training parameters used to produce the results using the microcircuit architecture.

| Parameter name                 |                      | Fig. 3                  | Fig. 5              |
|--------------------------------|----------------------|-------------------------|---------------------|
| <b>Neuron parameters</b>       |                      |                         |                     |
| $C_m$                          |                      | 1                       | 1                   |
| $E_l$                          |                      | 0                       | 0                   |
| $g_l$                          | $[\text{ms}^{-1}]^1$ | 0.03                    | 0.03                |
| $g^{\text{bas}}$               | $[\text{ms}^{-1}]$   | 0.1                     | 0.1                 |
| $g^{\text{api}}$               | $[\text{ms}^{-1}]$   | 0.06                    | 0.06                |
| $g^{\text{den}}$               | $[\text{ms}^{-1}]$   | 0.1                     | 0.1                 |
| $g^{\text{nudge, l}}$          | $[\text{ms}^{-1}]$   | 0.06                    | 0.06                |
| $g^{\text{nudge, tgt}}$        | $[\text{ms}^{-1}]$   | 0.06                    | 0.06                |
| $\tau_{\text{eff}}$            | $[\text{ms}]$        | $5.26^2$                | $5.26^2$            |
| $\varphi(x)$                   |                      | $\log[1 + \exp(x)]$     | $\log[1 + \exp(x)]$ |
| <b>Network parameters</b>      |                      |                         |                     |
| size input                     |                      | 9                       | 1                   |
| size hidden layer              |                      | 30                      | 1                   |
| size output layer              |                      | 3                       | 1                   |
| selfpred. $\eta_1^{\text{PP}}$ | $[\text{ms}^{-1}]$   | —                       | 0                   |
| selfpred. $\eta_2^{\text{PP}}$ | $[\text{ms}^{-1}]$   | —                       | 0                   |
| selfpred. $\eta_1^{\text{IP}}$ | $[\text{ms}^{-1}]$   | —                       | 40                  |
| selfpred. $\eta_1^{\text{PI}}$ | $[\text{ms}^{-1}]$   | —                       | 50                  |
| training $\eta_1^{\text{PP}}$  | $[\text{ms}^{-1}]$   | $5dt/T_{\text{pres}}^3$ | 50                  |
| training $\eta_2^{\text{PP}}$  | $[\text{ms}^{-1}]$   | $1dt/T_{\text{pres}}^3$ | 10                  |
| training $\eta_1^{\text{IP}}$  | $[\text{ms}^{-1}]$   | $2dt/T_{\text{pres}}^3$ | 20                  |
| training $\eta_1^{\text{PI}}$  | $[\text{ms}^{-1}]$   | $0^4$                   | $0^4$               |
| weight init (uniform)          |                      | $[-1, 1]$               | $[-1, 1]$           |
| <b>Training parameters</b>     |                      |                         |                     |
| start in selfpred. state       |                      | yes                     | no                  |
| train biases                   |                      | no                      | no                  |
| delay on target signal         |                      | $1dt^5$                 | $1dt^5$             |
| selfpred. epochs               |                      | —                       | 3                   |
| training epochs                |                      | 1000                    | 500                 |
| $dt$                           | $[\text{ms}]$        | 0.1                     | 0.01                |
| $T_{\text{pres}}$              |                      | $3dt - 5000dt$          | $100dt$             |

<sup>1</sup> To keep the other variables unitless, except for dynamical time scales, conductances and learning rates need to have the unit  $1/\text{ms}$ .

<sup>2</sup> The effective time constant is calculated from the neurons conductances:  $\tau_{\text{eff}} = \frac{C_m}{g_l + g^{\text{bas}} + g^{\text{api}}}$ .

<sup>3</sup> Learning rates are scaled with varying  $T_{\text{pres}}$ .

<sup>4</sup> If the network is starting in or has previously learned the self-predicting state the weights  $\mathbf{W}^{\text{PI}}$  do not need to be adapted.

<sup>5</sup> As the effect of a change in the input signal needs as many timesteps as there are hidden layers to reach the top layer of the network, the target signal needs to be delayed relative to the input by this amount of time steps.

## 15.3 Parameters used for the experiments shown in Fig. 4

Parameters not mentioned here explicitly (batch size, number of training, validation and test samples, learning rates, mean and variance for initial weights and biases) were the same as for the LE experiments shown in Fig. 2b (cf. Table 1).



**Table 3:** Additional parameters needed to reproduce the experiments shown in Fig. 4

| Symbol                     | Parameter name                   | Fig. 4a   | Fig. 4c                       |
|----------------------------|----------------------------------|---|-------------------------------|
| <b>Training parameters</b> |                                  |   |                               |
| $dt$                       | temporal resolution              | 0.002 ms – 2.0 ms                               | 0.01 ms, 0.05 ms <sup>1</sup> |
| $T_{\text{pres}}$          | presentation time                | 100 $dt$  | 20 $dt$ – 1000 $dt$           |
| $\tau^s$                   | synaptic time constant           | –   | 0 ms – 2.0ms                  |
|                            | # seeds                          | 4   | 3                             |
| <b>Noise parameters</b>    |                                  |   |                               |
| $\sigma_{\tau^{m/r}}$      | time constant width <sup>2</sup> | $0 \tau^{m/r}, 0.01 \tau^{m/r}, 0.2 \tau^{m/r}$ | –                             |
| $\sigma_{\xi}$             | noise width                      | –   | $0.2 r_{\text{max}}$          |
|                            | target LPF                       | –   | yes <sup>3</sup>              |

<sup>1</sup> The smaller value was used for the simulations of the green datapoints while the bigger value was used to obtain the blue and yellow curves.

<sup>2</sup> Time constants were clipped, i.e.,  $\tau^{m/r} + \xi \in [1, 1000]$ , to exclude the unphysical case of them to become negative.

<sup>3</sup> In case of Fig. 4c, an additional LPF with time constant  $\tau^{\text{tg}} = N \times \tau^s$  where  $N = \#$  layers was applied to the target signal. However, this is just an approximation for the  $N$  LPFs that are being applied to the input signal during a forward pass. Yet it helps to reduce “wrong learning” during the short relaxation phases introduced by the synaptic filtering and can be neglected in the limit  $\tau^s \rightarrow 0$ .

## 16 Broader impact

The physical interpretation of our model not only offers a biologically plausible implementation of continuous-time, continuously active neuro-synaptic dynamics, but also outlines a specific path towards mixed-signal (analog/digital) in-silico implementation. Even with existing technologies, such neuromorphic systems harbor the potential of surpassing their biological archetypes with respect to both energy efficiency and speed [2]. In conjunction with the inherent ability of our framework to support the processing of continuous data streams, the reduced power consumption of such devices makes them a prime candidate for the construction of autonomous, embodied learning machines.

While obviously beneficial for research and commercial deployment, one should be aware that improved training efficiency carries the risk of deploying ever more intransparent models [3]. Furthermore, along with its obvious benefits, improved, and in particular autonomous AI entails a plethora of far-reaching societal consequences that are the subject of ongoing academic and public debate [4]. Future progress hence needs to be considered carefully and responsibly, and, in particular, properly reflected in public policy.

On the path towards understanding and replicating biological intelligence, a corollary benefit for the scientific community may emerge. Modern machine learning requires enormous amounts of compute, thus largely limiting cutting-edge developments to institutions with the corresponding resources. The envisioned bio-inspired yet also bio-transcendent hardware systems have the potential to drastically increase the overall efficiency of custom-designed computational platforms. The resulting decrease in operating costs could thus significantly expedite the democratization of AI research.

## References

1. Sacramento, J., Ponte Costa, R., Bengio, Y. & Senn, W. Dendritic Cortical Microcircuits Approximate the Backpropagation Algorithm. *Advances in Neural Information Processing Systems* **31**, 8721–8732 (2018).
2. Marković, D., Mizrahi, A., Querlioz, D. & Grollier, J. Physics for Neuromorphic Computing. *Nature Reviews Physics* **2**, 499–510 (2020).
3. Bender, E. M., Gebru, T., McMillan-Major, A. & Shmitchell, S. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 610–623 (2021).

4. Futurium (European Commission). *Ethics Guidelines for Trustworthy AI* <https://ec.europa.eu/futurium/en/ai-alliance-consultation>.