# Dendritic error backpropagation
# in deep cortical microcircuits

João Sacramento[1*], Rui Ponte Costa[1], Yoshua Bengio[2], Walter Senn[1*]

[1]Department of Physiology
University of Bern, Switzerland

[2]Montreal Institute for Learning Algorithms
Université de Montréal, Quebec, Canada

## Abstract

Animal behaviour depends on learning to associate sensory stimuli with the desired motor command. Understanding how the brain orchestrates the necessary synaptic modifications across different brain areas has remained a longstanding puzzle. Here, we introduce a multi-area neuronal network model in which synaptic plasticity continuously adapts the network towards a global desired output. In this model synaptic learning is driven by a local dendritic prediction error that arises from a failure to predict the top-down input given the bottom-up activities. Such errors occur at apical dendrites of pyramidal neurons where both long-range excitatory feedback and local inhibitory predictions are integrated. When local inhibition fails to match excitatory feedback an error occurs which triggers plasticity at bottom-up synapses at basal dendrites of the same pyramidal neurons. We demonstrate the learning capabilities of the model in a number of tasks and show that it approximates the classical error backpropagation algorithm. Finally, complementing this cortical circuit with a disinhibitory mechanism enables attention-like stimulus denoising and generation. Our framework makes several experimental predictions on the function of dendritic integration and cortical microcircuits, is consistent with recent observations of cross-area learning, and suggests a biological implementation of deep learning.

## Introduction

While walking on the street we are constantly bombarded with complex sensory stimuli. Learning to navigate such complex environments is of fundamental importance for survival. In the brain, these forms of learning are believed to rely on the orchestrated wiring of synaptic communication between different cortical areas, such as visual and motor cortices (Petreanu et al., 2012; Manita et al., 2015; Makino and Komiyama, 2015; Poort et al., 2015; Zmarz and Keller, 2016; Attinger et al., 2017). However, how to correctly modify synapses to achieve an appropriate interaction between brain areas has remained an open question. This fundamental issue in learning and development is often referred to as the credit assignment problem (Rumelhart et al., 1986; Sutton and Barto, 1998; Roelfsema and van Ooyen, 2005; Friedrich et al., 2011; Bengio, 2014). The brain,

---
* Corresponding authors: {sacramento},{senn}@pyl.unibe.ch

and artificial neural networks alike, have to determine how to best modify a given synapse across multiple processing stages to ultimately improve global behavioural output.

Machine learning has recently undergone remarkable progress through the use of deep neural networks, leading to human-level performance in a growing number of challenging problems (LeCun et al., 2015). Key to an overwhelming majority of these achievements has been the backpropagation of errors algorithm (backprop; Rumelhart et al., 1986), which has been long dismissed in neuroscience on the grounds of biologically implausibility (Grossberg, 1987; Crick, 1989). Nonetheless, accumulating evidence highlights the difficulties of simpler learning models and architectures in accurately reproducing cortical activity patterns when compared to deep neural networks, notably trained only on sensory data (Yamins et al., 2014; Khaligh-Razavi and Kriegeskorte, 2014; Yamins and DiCarlo, 2016). Although recent developments have started to bridge the gap between neuroscience and artificial intelligence (Marblestone et al., 2016; Lillicrap et al., 2016; Costa et al., 2017; Guerguiev et al., 2017), whether the brain implements a backprop-like algorithm remains unclear.

Here we propose that the errors at the heart of backprop are encoded on the distal dendrites of cross-area projecting pyramidal neurons. In our model, these errors arise from a failure to exactly match via lateral (e.g. somatostatin-expressing, SST) interneurons the top-down feedback from downstream cortical areas. Synaptic learning is driven by these error-like signals that flow through the dendrites and trigger plasticity on bottom-up connections. Therefore, in contrast to previous approaches (Marblestone et al., 2016), in our framework a given neuron is used simultaneously for activity propagation (at the somatic level), error encoding (at distal dendrites) and error propagation to the soma. Importantly, under certain simplifying assumptions, we were able to formally show that learning in the model approximates backprop.

We first illustrate the different components of the model and afterwards demonstrate its performance by training a multi-area network on associative nonlinear regression and recognition tasks (handwritten digit image recognition, a standard benchmark). Then, we further extend the framework to consider learning of the top-down synaptic pathway. When coupled with a disinhibitory mechanism this allows the network to generate prototypes of learnt images as well as perform input denoising. We interpret this disinhibitory mechanism as being implemented through another inhibitory cell-type (e.g. vasoactive intestinal peptide-expressing, VIP interneurons). Finally, we make several experimentally testable predictions in terms of the role of dendrites and different interneuron types being involved while an animal learns to associate signals originating from different brain areas.

# Results

The cortex exhibits remarkably intricate, but stereotypical circuits. Below we describe a plastic cortical circuit model that considers two features observed in neocortex: dendritic compartments and different cell types. Cross-area synapses onto the dendritic compartments learn to reduce the prediction error between the somatic potential and their own dendritic branch potential. Additionally, lateral synaptic input from local interneurons learns to cancel top-down feedback from downstream brain areas. When a new top-down input arrives at distal dendrites that cannot be matched by lateral inhibition it signals a neuron-specific error (encoded on the dendritic potential) that triggers synaptic learning at a given pyramidal cell. As learning progresses, the interneurons gradually learn to cancel once again the new input until eventually learning stops. We show that this cortical circuit implements error backpropagation, and demonstrate its performance in various tasks.

## The dendritic cortical circuit learns to predict self-generated top-down input

We first study a generic network model with $N$ cortical brain areas (a multi-layer network, in machine learning parlance), comprising an input area (representing, for instance, thalamic input to sensory areas, denoted as area 0), one or more 'hidden' areas (representing secondary sensory and consecutive higher brain areas, denoted as area $k$ and area $k+1$, respectively) and an output brain area (e.g. motor cortex, denoted as area $N$), see schematic in Fig. 1A. Unlike conventional artificial neural networks, hidden neurons feature both bottom-up ($\mathbf{W}_{k+1,k}^{PP}$) and top-down ($\mathbf{W}_{k,k+1}^{PP}$) connections, thus defining a recurrent network structure. Top-down synapses continuously feed back the next brain area predictions to a given bottom-up input. Our model uses of this feedback to determine corrective error signals and ultimately guide synaptic plasticity across multiple areas.

Building upon previous work (Urbanczik and Senn, 2014), we adopt a simplified multicompartment neuron and describe pyramidal neurons as three-compartment units (schematically depicted in Fig. 1A; see also Methods). These compartments represent the somatic, basal and apical integration zones that characteristically define neocortical pyramidal cells (Spruston, 2008; Larkum, 2013). The dendritic structure of the model is exploited by having bottom-up and top-down synapses converging onto separate dendritic compartments (basal and distal dendrites, respectively), consistent with experimental observations (Spruston, 2008) and reflecting the preferred connectivity pattern of cortico-cortical projections (Larkum, 2013).

Consistent with the neurophysiology of SST interneurons (Urban-Ciecko and Barth, 2016), we
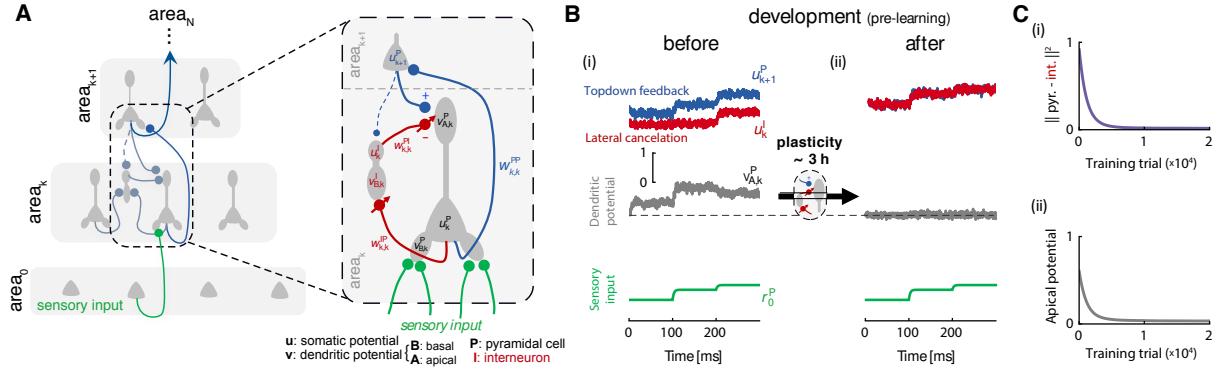
Figure 1: **Dendritic cortical circuit learns to predict self-generated top-down input.** (**A**) Illustration of multi-area network architecture. The network consists of an input area 0 (e.g., thalamic input), one or more intermediate (hidden) areas (represented by area $k$ and area $k+1$, which can be mapped onto primary and higher sensory areas) and an output area $N$ (e.g. motor cortex) (left). Each hidden area consists of a microcircuit with pyramidal cells and lateral inhibitory interneurons (e.g. somatostatin-positive, SST, cells) (right). Pyramidal cells consist of three compartments: a basal compartment (with voltage $\mathbf{v}_{B,k}^{P}$) that receives bottom-up input; an apical compartment ($\mathbf{v}_{A,k}^{P}$), where top-down input converges to; and a somatic compartment ($\mathbf{u}_{k}^{P}$), that integrates the basal and apical voltage. Interneurons receive input from lateral pyramidal cells onto their own basal dendrites ($\mathbf{v}_{B,k}^{I}$), integrate this input on their soma ($\mathbf{u}_{k}^{I}$) and project back to the apical compartments ($\mathbf{v}_{A,k}^{P}$) of same-area pyramidal cells. (**B**) In a pre-learning developmental stage, the network learns to predict and cancel top-down feedback given randomly generated inputs. Only pyramidal-to-interneuron synapses ($\mathbf{W}_{k,k}^{IP}$) and interneuron-to-pyramidal synapses ($\mathbf{W}_{k,k}^{PI}$) are changed at that stage according to predictive synaptic plasticity rules (defined in Eqs 8 and 9 of the Methods). Example voltage traces for a randomly chosen downstream neuron ($\mathbf{u}_{k+1}^{P}$) and a corresponding interneuron ($\mathbf{u}_{k}^{I}$), a pyramidal cell apical compartment ($\mathbf{v}_{A,k}^{P}$) and an input neuron ($\mathbf{u}_{0}^{P}$), before (i) and after (ii) development, for three consecutively presented input patterns. Once learning of the lateral synapses from and onto interneurons has converged, self-generated top-down signals are predicted by the network — it is in a *self-predicting state*. Here we use a concrete network with one hidden area and 30-20-10 pyramidal neurons (input-hidden-output), but the particular network dimensions do not impact the ability of the network to produce these results. Note that no desired targets are presented to the output area (cf. Fig. 2); the network is solely driven by random inputs. (**C**) Lateral inhibition cancels top-down input. (i) Interneurons learn to match next-area pyramidal neuron activity as their input weights $\mathbf{W}_{k,k}^{IP}$ adapt (see main text and Methods for details). (ii) Concurrently, learning of interneuron-to-pyramidal synapses ($\mathbf{W}_{k,k}^{PI}$) silences the apical compartment of pyramidal neurons, but pyramidal neurons remain active (cf. B). This is a general effect, as the lateral microcircuit learns to predict and cancel the expected top-down input for every random pattern (see SI).

also introduce a second population of cells within each hidden area with both lateral and cross-area connectivity, whose role is to cancel the top-down input. Modelled as two-compartment units (depicted in red, Fig. 1A; see also Methods), such interneurons are predominantly driven by pyramidal cells within the same area through weights $\mathbf{W}_{i,i}^{IP}$, and they project back to the apical dendrites of the same-area pyramidal cells through weights $\mathbf{W}_{k,k}^{PI}$ (see Fig. 1A). Additionally, cross-area feedback onto SST cells originating at the next higher brain area $k+1$ provide a weak nudging signal for these interneurons, modelled after Urbanczik and Senn (2014) as a conductance-based so-

matic input current. For computational simplicity, we modelled this weak top-down nudging on a one-to-one basis (that can also be relaxed): each interneuron is nudged towards the potential of a corresponding upper-area pyramidal cell. Recent monosynaptic input mapping experiments show that somatostatin-positive cells (SST, of which Martinotti cells are the main type) in fact receive also top-down projections (Leinweber et al., 2017), that according to our proposal encode the weak 'teaching' signals from higher to lower brain areas.

As detailed below, this microcircuit is key to encode and backpropagate errors across the network. We first show how synaptic plasticity of lateral interneuron connections establishes a network regime, which we term *self-predicting*, whereby lateral input cancels the self-generated top-down feedback, effectively silencing apical dendrites. For this reason, SST cells are functionally inhibitory and are henceforth referred to as interneurons. Crucially, when the circuit is in this so-called self-predicting state, presenting a novel external signal at the output area gives rise to top-down activity that cannot be explained away by the interneuron circuit. Below we show that these apical mismatches between top-down and lateral input constitute the backpropagated, neuron-specific errors that drive plasticity on the forward weights to the hidden pyramidal neurons.

Learning to predict the feedback signals involves adapting both weights from and to the lateral interneuron circuit. Consider a network that is driven by a succession of sensory input patterns (Fig. 1B, bottom row). The exact distribution of inputs is unimportant as long as they span the whole input space (see SI). Learning to cancel the feedback input is divided between both the weights from pyramidal cells to interneurons, $\mathbf{W}_{k,k}^{\mathrm{IP}}$, and from interneurons to pyramidal cells, $\mathbf{W}_{k,k}^{\mathrm{PI}}$.

First, due to the somatic teaching feedback, learning of the $\mathbf{W}_{k,k}^{\mathrm{IP}}$ weights leads interneurons to better reproduce the activity of the respective higher brain area $k+1$ (Fig. 1B (i)). A failure to reproduce area $k+1$ activity generates an internal prediction error at the dendrites of the interneurons, which triggers synaptic plasticity (as defined by Eq. 8 in the Methods) that corrects for the wrong dendritic prediction and eventually leads to a faithful tracing of the upper area activity by the lower area interneurons (Fig. 1B (ii)). The mathematical analysis (see SI, Eq. S27) shows that the plasticity rule (8) makes the inhibitory population implement the same function of the area-$k$ pyramidal cell activity as done by the area–$(k+1)$ pyramidal neurons. Thus, the interneurons will learn to mimic the area–$(k+1)$ pyramidal neurons (Fig. 1Ci).

Second, as the interneurons mirror upper area activity, inter-to-pyramidal neuron synapses within the same area ($\mathbf{W}_{k,k}^{\mathrm{PI}}$, Eq. 9) successfully learn to cancel the top-down input to the apical dendrite (Fig. 1Cii), independently of the actual input stimulus that drives the network. By doing

so, the inter-to-pyramidal neuron weights $\mathbf{W}_{k,k}^{\text{PI}}$ learn to mirror the top-down weights onto the lower area pyramidal neurons. The learning of the weights onto and from the interneurons works in parallel: as the interneurons begin to predict the activity of pyramidal cells in area $k+1$, it becomes possible for the plasticity at interneuron-to-pyramidal synapses (Eq. 9) to find a synaptic weight configuration which precisely cancels the top-down feedback (see also SI, Eq. S29). At this stage, every pattern of activity generated by the hidden areas of the network is explained by the lateral circuitry, Fig. 1C (ii). Importantly, once learning of the lateral interneurons has converged, the apical input cancellation occurs irrespective of the actual bottom-up sensory input. Therefore, interneuron synaptic plasticity leads the network to a *self-predicting state*. We propose that the emergence of this state could occur during development, consistent with experimental findings (Dorrn et al., 2010; Froemke, 2015). Starting from a cross-area self-predicting configuration helps learning of specific tasks (but is not essential, see below and Methods).

**Deviations from self-predictions encode backpropagating errors**

Having established a self-predicting network, we next show how prediction errors get propagated backwards when a novel input is provided to the output area. This new signal, which we model via the activation of additional somatic conductances in output pyramidal neurons (see Methods), plays the role of a teaching or associative signal (see specific tasks below). Here we consider a concrete implementation of the network model introduced above, with an input, a hidden and an output brain area (areas 0, 1 and 2, respectively; Fig. 2A). We demonstrate learning in the model with a simple task: memorizing a single input-output pattern association. This setup naturally generalizes to multiple memories by iterating over a set of associations to be learned.

When the pyramidal cell activity in the output area is nudged towards some desired target (Fig. 2B (i)), the bottom-up synapses $\mathbf{W}_{2,1}^{\text{PP}}$ from the lower area neurons to the basal dendrites are adapted, again according to the plasticity rule that implements the dendritic prediction of somatic spiking (see Eq. 7 in the Methods and Urbanczik and Senn (2014)). What these synapses cannot explain away shows up as a dendritic error in the pyramidal neurons of the lower area 1. In fact, the self-predicting microcircuit can only cancel the feedback that is produced by the lower area activity. Due to the unexplained teaching signal in the output area, the top-down input partially survives the lateral inhibition; this leads to the activation of distal dendrites (Fig. 2B (i)). The mathematical analysis reveals that the apical deviation from baseline encodes an error that is effectively backpropagated from the output area.
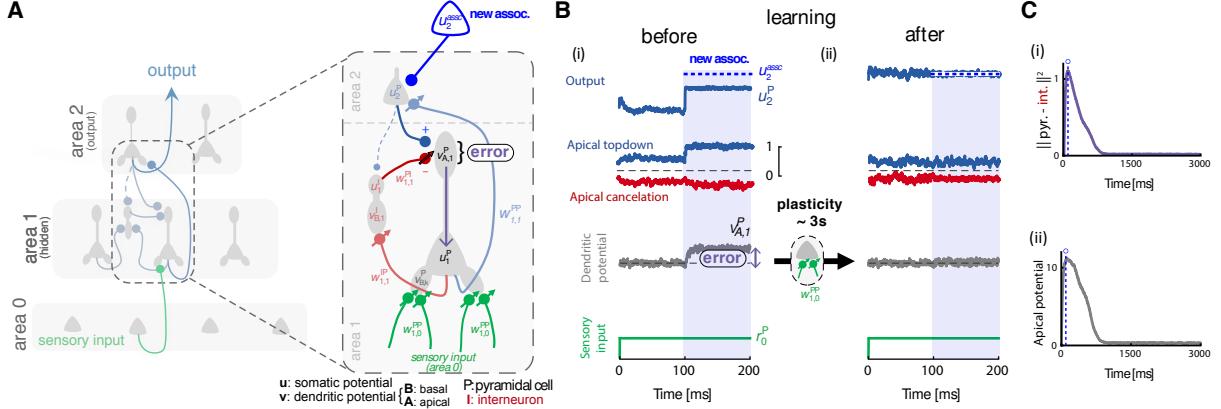
6

Figure 2: **Deviations from self-predictions encode backpropagating errors that are used for learning in bottom-up synapses.** (**A**) When a novel associative (or 'teaching') signal is presented to the output area ($\mathbf{u}_2^{\text{assoc}}$, blue at the top), a prediction error in the apical compartments of pyramidal neurons in the upstream area (area 1, 'error') is generated. This error appears as an apical voltage deflection that propagates down to the soma (purple arrow) where it modulates the somatic firing rate. Bottom-up synapses at the basal dendrites learn to predict the somatic firing rate (bottom, green). Only the elements directly involved in encoding the error and modifying the bottom-up synapses are highlighted in this microcircuit. (**B**) Activity traces in the microcircuit before and after a new associative signal is learned. (i) Before learning, when a new associative signal is presented at the output area ($\mathbf{u}_2^{\text{assoc}}$, blue dashed), the somatic voltage of output neurons changes accordingly ($\mathbf{u}_2^{\text{P}}$, grey blue, top). As a result, a new signal is observed in the upstream apical dendrite ($\mathbf{v}_{A,1}^{\text{P}}$, grey bottom) which originates from a mismatch between the top-down feedback (grey blue) and the cancellation given by the lateral interneurons (red). (ii) After $\sim$3s of learning with this new associative signal, plasticity at the bottom-up synapses ($\mathbf{W}_{1,0}^{\text{PP}}$), which receive sensory input ($r_0^{\text{P}}$, green), leads to a near-exact prediction of the new, previously unexplained associative signal $\mathbf{u}_2^{\text{assoc}}$ by $\mathbf{u}_2^{\text{P}}$. Consequently, the distal dendrite no longer shows a voltage deflection (error), which results from the top-down and lateral cancellation inputs having the magnitude but opposite signs (blue and red traces, respectively). The network now fully reproduces the new associative signal (top). (**C**) Learning gradually explains away the backpropagated activity. (i) Interneurons learn to predict, and effectively cancel, the backpropagated activity as the lateral weights from the pyramidal-to-interneurons ($\mathbf{W}_{1,1}^{\text{IP}}$) adapt. (ii) While simultaneously adapting the interneuron-to-pyramidal synapses ($\mathbf{W}_{1,1}^{\text{PI}}$), the apical compartment is eventually silenced, even though the pyramidal neurons remain active (not shown). Vertical blue dashed line represents the moment when the new associative signal is presented for the first time.

The somatic integration of apical activity induces plasticity at the bottom-up synapses $\mathbf{W}_{1,0}^{\text{PP}}$ on the basal dendrites. As described above, plasticity at these synapses too is governed by the dendritic prediction of somatic activity, just as for the synapses to the interneurons (Eq. 7). As the apical error changes the somatic activity, plasticity of the $\mathbf{W}_{1,0}^{\text{PP}}$ weights tries to further reduce the error in the output area. Importantly, the plasticity rule depends only on information that is available at the synaptic level. More specifically, it is a function of both postsynaptic firing and dendritic branch voltage, as well as presynaptic activity, in par with detailed phenomenological models (Clopath et al., 2010; Bono and Clopath, 2017). In a spiking neuron model, the plasticity rule can reproduce a

number of experimental results on spike-timing dependent plasticity (Spicher et al., in preparation).

In contrast to the establishing of the self-predicting network state, learning now involves the simultaneous modifications of both lateral circuit and bottom-up synaptic weights (Fig. 2). On the one hand, lateral weights track changes in output area activity, in this way approximately maintaining the network in a self-predicting state throughout learning. On the other hand, the inputs to area 1 pyramidal neurons adapt to reduce prediction errors. Altogether, plasticity eventually leads to a network configuration in which the novel top-down input is successfully predicted (Fig. 2B,C).

## Cross-area network learns to solve a nonlinear associative task

So far we have described the key components of our model in a multi-area network using a toy problem. Now, we turn to more challenging problems. The first of which is a nonlinear associative task, where the network has to learn to associate the sensory input with the output of a separate multi-area network that transforms the same sensory input — this can be recast as a nonlinear regression problem (Fig. 3A; see Methods for details on the architecture and learning conditions).
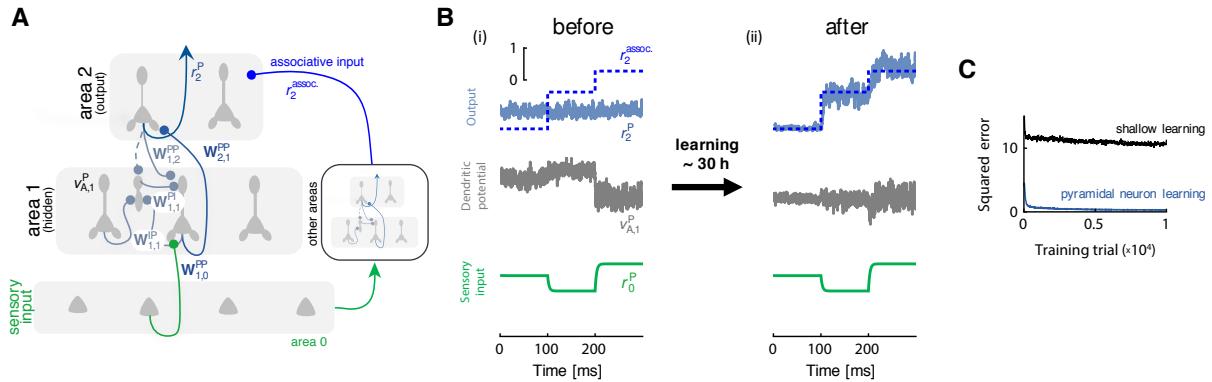


Figure 3: **A multi-area network learns to solve a nonlinear associative task online in continuous time and without phases.** (**A**) Starting from a self-predicting network state (cf. Fig. 1), a 30-20-10 fully-connected pyramidal neuron network learns to approximate a nonlinear function (represented by the 'other areas' box, another 30-20-10 network) from a stream of input-output pattern pairs. The neuronal and synaptic weight dynamics evolve continuously and without interruption. (**B**) Example firing rates for a randomly chosen output neuron ($r_2^P$, blue noisy trace) and its desired target imposed by the associative input ($r_2^{assoc}$, blue dashed line), together with the voltage in the apical compartment of a hidden neuron ($v_{A,1}^P$, grey noisy trace) and the input rate from the sensory neuron ($r_0^P$, green). Before learning (i), the apical dendrite of the hidden neuron shows errors in response to three consecutive input patterns that disappear after 30 h of successful learning (ii). The presentation of the novel output target produces deviations from baseline in the apical compartment, visible in the initial $v_{A,1}^P$ trace. Such mismatch signals trigger plasticity in the forward weights $\mathbf{W}_{1,0}^{PP}$, which eventually leads to a reduction of the error in the output area, and henceforth a return to baseline of the apical voltage in the hidden area below. (**C**) Error curves for the full model and a shallow learner for comparison, where no backpropagation of errors occurs and only the output weights $\mathbf{W}_{2,1}^{PP}$ are adapted.

8

We let learning occur in continuous time without pauses or alternations in plasticity as input patterns are sequentially presented. This is in contrast to previous learning models that rely on computing activity differences over distinct phases, requiring temporally nonlocal computation, or globally coordinated plasticity rule switches (Hinton and McClelland, 1988; O'Reilly, 1996; Xie and Seung, 2003; Scellier and Bengio, 2017; Guerguiev et al., 2017). Furthermore, we relaxed the bottom-up vs. top-down weight symmetry imposed by the backprop algorithm and kept the top-down weights $\mathbf{W}_{1,2}^{\mathrm{PP}}$ fixed. Feedback $\mathbf{W}_{1,2}^{\mathrm{PP}}$ weights quickly aligned to $\sim 45^{\mathrm{o}}$ of the forward weights $\left(\mathbf{W}_{2,1}^{\mathrm{PP}}\right)^T$, in line with the recently discovered feedback alignment phenomenon (Lillicrap et al., 2016). This simplifies the architecture, because top-down and interneuron-to-pyramidal synapses need not to be changed. Finally, to test the robustness of the network, we injected a weak noise current to every neuron, as a simple model for uncorrelated background activity (see Methods). Our network was still able to learn this harder task (Fig. 3B), performing considerably better than a shallow learner where only output weights were adjusted (Fig. 3C). Useful changes were thus made to hidden area 1 bottom-up weights; the network effectively solved the credit assignment problem.

## Multi-area network learns to discriminate handwritten digits

Next, we turn to a standard machine learning problem, the classification of handwritten digits from the MNIST database. This data set is popularly used to study the performance of learning models, including various artificial neural networks trained with backprop. Notably, shallow models (e.g., logistic regression) or networks trained with plain Hebbian learning alone suffer from poor performance and do not offer a remedy for the problem.

We wondered how our model would fare in this real-world benchmark, in particular whether the prediction errors computed by the interneuron microcircuit would allow learning the weights of a hierarchical nonlinear network with multiple hidden areas. To that end, we trained a deeper, larger 4-area network (with 784-500-500-10 pyramidal neurons, Fig. 4A) by pairing digit images with teaching inputs that nudged the 10 output neurons towards the correct class pattern. To speed up the experiments we studied a simplified network dynamics which determined compartmental potentials without requiring a full neuronal relaxation procedure (see Methods). As in the previous experiments, synaptic weights were randomly initialized and set to a self-predicting configuration where interneurons cancelled top-down inputs, rendering the apical compartments silent before training started. Top-down and interneuron-to-pyramidal weights were kept fixed.

The network was able to achieve a test error of 1.96%, Fig. 4B, a figure not overly far from the
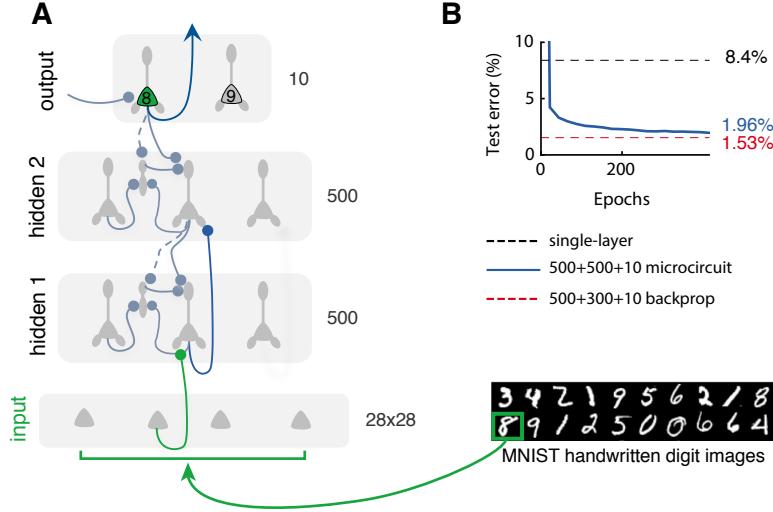
Figure 4: **Learning to classify real-world, structured stimuli with a multi-area network.** (**A**) A 784-500-500-10 (i.e. with two hidden areas) network of pyramidal neurons learns to recognize and classify handwritten digits from the MNIST data set. Only a subset of connections is shown to enhance clarity. (**B**) Competitive accuracy ($< 2\%$, an empirical signature of backprop-like learning) is achieved on the standard MNIST testing dataset by our network (solid blue). For comparison the performance of a shallow learner (i.e. a network in which only output weights are adapted, dashed black) and of a standard artificial neural network trained with backprop (dashed red, see Methods) are also shown.

reference mark of non-convolutional artificial neural networks optimized with backprop (1.53%) and comparable to recently published results that lie within the range 1.6-2.4% (Lee et al., 2015; Lillicrap et al., 2016). This was possible even though interneurons had to keep track of changes to forward weights as they evolved, simultaneously and without phases. Indeed, apical compartment voltages remained approximately silent when output nudging was turned off (data not shown), reflecting the maintenance of a self-predicting state throughout learning. Moreover, thanks to a feedback alignment dynamics (Lillicrap et al., 2016), the interneuron microcircuit was able to translate the feedback from downstream areas into single neuron prediction error signals, despite the asymmetry of forward and top-down weights and at odds with exact backprop.

**Disinhibition enables sensory input generation and sharpening**

So far we assumed that feedback from downstream neurons is relayed through fixed top-down synapses. However, this need not be so. As we demonstrate next, the interneuron microcircuit is capable of tracking changes to the top-down stream dynamically as learning progresses. This endows the model with important additional flexibility, as feedback connections — known to mediate attention and perceptual acuity enhancement in sensory cortices — are likely plastic (Huber et al.,

2012; Petreanu et al., 2012; Manita et al., 2015; Makino and Komiyama, 2015; Attinger et al., 2017; Leinweber et al., 2017).
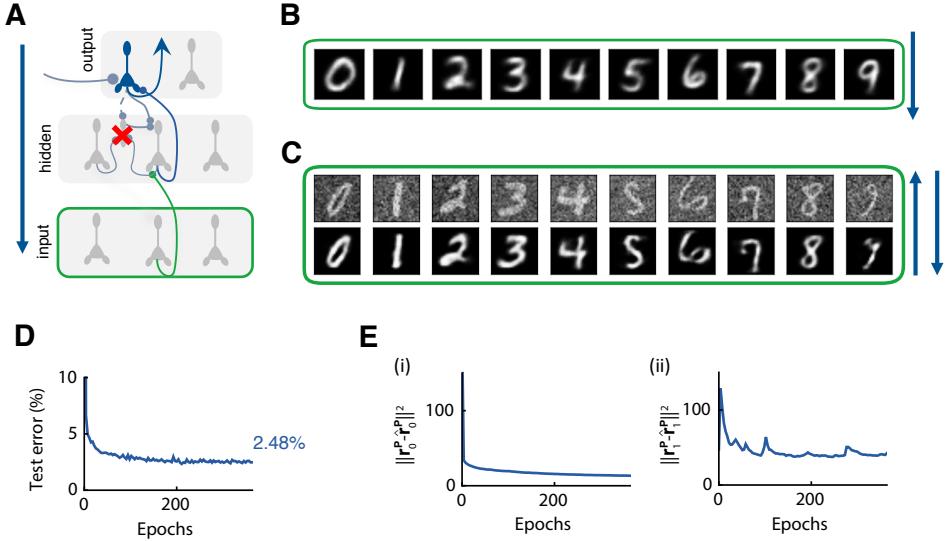


Figure 5: **Top-down synapses can be adapted to simultaneously drive bottom-up learning, input construction and denoising.** (**A**) Classification performance of a 784-1000-10 network exposed to MNIST images, with plastic top-down synapses that learns to predict lower-area activities. Top-down and forward weights co-adapt without pauses or phases. (**B**) Driving the network top-to-bottom (i.e., initializing the output area to a particular digit and turning off lateral and bottom-up inputs of both hidden and input areas) recreates class-specific image examples in the input area. The top-down connections can be tuned to encode a simple inverse visual model. (**C**) Such an inverse model yields image denoising, which can be achieved by reconstructing corrupted inputs from hidden area activities. (**D**) The network also successfully learns to classify images. (**E**) Inverse reconstruction losses of original images (i) and hidden (ii) neuron activities. Top-down synapses connecting hidden pyramidal neurons back to the input area learn to reconstruct pixel arrangements given hidden neuron activities; synapses originating at the output area learn to predict hidden area activities given the current class label estimate.

As a case in point we considered a simple extension to a three-area network of 784-1000-10 pyramidal neurons again exposed to MNIST images, Fig. 5. The architecture is as before, except that we now let dendritic predictive plasticity shape the top-down weights from output to hidden neurons $\mathbf{W}_{1,2}^{\mathrm{PP}}$ as well as an extra set of weights $\mathbf{W}_{0,1}^{\mathrm{PP}}$ connecting hidden neurons back to the input area (see Eq. 10 in the Methods).

In this extended network, top-down synapses learn to predict the activities of the corresponding area below and thus implement an approximate inverse of the forward model. In effect, these connections play a dual role, beyond their sole purpose in backprop. They communicate back upper area activities to improve the hidden neuron representation on a recognition task, and they learn to invert the induced forward model. This paired encoder-decoder architecture is known as target propagation in machine learning (Bengio, 2014; Lee et al., 2015). Our compartmental pyramidal

11

neuron model affords a simple design for the inverse learning paradigm: once more, plasticity of top-down synapses is driven by a postsynaptic dendritic error factor, comparing somatic firing with a local branch potential carrying the current top-down estimate.

Importantly, our results show that the network still learned to recognize handwritten digits, Fig. 5A, reaching a classification error of 2.48% on the test set. This again highlights that not only transposed forward weight matrices, as prescribed by backprop, deliver useful error signals to hidden areas. In this experiment, we initialized every weight matrix randomly and independently, and did not pre-learn lateral circuit weights. Although forward, top-down and lateral weights were all jointly adapted starting from random initial conditions, a self-predicting state quickly ensued, leading to a drop in classification error. Concomitantly, the reconstructions of hidden neuron activities and input images improved, Fig. 5B.

The learned inverse model can be used to generate prototypical digit images in the input area. We examined qualitatively its performance by directly inspecting the produced images. Specifically, for each digit class we performed a top-to-bottom pass with lateral inhibition turned off, starting from the corresponding class pattern $\mathbf{r}_2^P$. For simplicity, we disabled basal feedforward inputs as well to avoid recurrent effects (see Methods). This procedure yielded prototype reconstructions $\hat{\mathbf{r}}_0^P = \phi(\mathbf{W}_{0,1}^{PP} \phi(\mathbf{W}_{1,2}^{PP} \mathbf{r}_2^P))$ which resemble natural handwritten digits, Fig. 5C, confirming the observed decrease in reconstruction loss.

Crucially, for the network to be able to generate images, the apical dendrites of hidden neurons should be fully driven by their top-down inputs. In terms of our microcircuit implementation, this is achieved by momentarily disabling the contributions from lateral interneurons. A switch-like dis-inhibition (Pi et al., 2013) is thus capable of turning apical dendrites from error signalling devices into regular prediction units: the generative mode corresponds to a disinhibited mode of operation. Due to their preferential targetting of SST interneurons, VIP interneurons are likely candidates to implement this switch.

Recent reports support the view that cortico-cortical feedback to distal dendrites plays an active role as mice engage in perceptual discrimination tasks (Manita et al., 2015; Makino and Komiyama, 2015; Takahashi et al., 2016). Inspired by these findings, we further tested the capabilities of the model in a visual denoising task, where the prior knowledge incorporated in the top-down network weights is leveraged to improve perception. In Fig. 5D, we show the reconstructions $\hat{\mathbf{r}}_0^P = \phi(\mathbf{W}_{0,1}^{PP} \mathbf{r}_1^P)$ obtained after presenting randomly picked MNIST images from the test set that had been corrupted with additive Gaussian noise. We show only the apical predictions carried by top-down inputs back

to sensory area 0, without actually changing area 0 activity. Interestingly, we found that the hidden neuron representations shaped by classification errors served as reasonable visual features for the inverse model as well. Most of the noise was filtered out, although some of the finer details of the original images were lost in the process.

## Discussion

How the brain successfully assigns credit and modifies specific synaptic connections given a global associative signal has puzzled neuroscientists for decades. Here we have introduced a novel framework in which a single neuron is capable of transmitting predictions as well as prediction errors. These neuron-specific errors are encoded at distal dendrites and are kept in check by lateral (e.g. somatostatin-expressing) interneurons. Next, local synaptic plasticity mechanisms use such dendritic-encoded prediction errors to correctly adjust synapses. We have shown that these simple principles allow networks of multiple areas to successfully adjust their weights in challenging tasks, and that this form of learning approximates the well known backpropagation of errors algorithm.

**Experimental predictions**    Because our model touches on a number of levels: from brain areas to microcircuits, dendritic compartments and synapses, it makes several predictions. Here we highlight some of these predictions and related experimental observations:

(1) *Dendritic error representation.* Probably the most fundamental feature of the model is that dendrites, in particular distal dendrites, encode error signals that instruct learning of lateral and downstream connections. This means that during a task that required the association of two brain areas to develop, lateral interneurons would modify their synaptic weights such that the top-down signals are cancelled. Moreover, during learning, or if this association is broken, a dendritic error signal should be observed. While monitoring such dendritic signals during learning is challenging, there is recent experimental evidence that supports this model. Mice were trained in a simple visuomotor task where the visual flow in a virtual environment presented to the animal was coupled to its own movement (Zmarz and Keller, 2016; Attinger et al., 2017). When this coupling was broken (by stopping the visual flow) mismatch signals were observed in pyramidal cells, consistent with the prediction error signals predicted by our model.

(2) *Lateral inhibition of apical activity.* Our apical error representation is based on lateral inhibitory feedback to distal dendritic compartments of pyramidal cells. There is evidence for top-down

feedback to target distal (layer-1) synapses of both layer-2/3 and layer-5 pyramidal cells (Petreanu et al., 2009), and both cell types have lateral somatostatin interneurons which target the distal dendrites of the respective pyramidal cells (Markram et al., 2004). The cancellation of the feedback provided by somatostatin interneurons should be near-exact both in its magnitude and delay. In the brain, there can be a substantial delay between the lateral excitatory input and the feedback from other brain areas (in the order of tens to hundreds of milliseconds (Cauller and Kulics, 1991; Larkum, 2013)), suggesting that the lateral inhibitory interaction mediated by SST cells should be also delayed and tuned to the feedback. Interestingly, there is strong experimental support for a delayed inhibition mediated by pyramidal-to-SST connections (Silberberg and Markram, 2007; Murayama et al., 2009; Berger et al., 2009; Berger et al., 2010), which could in principle be tuned to match both the delay and magnitude of the feedback signal. Moreover, the spontaneous activity of SST interneurons is relatively high (Urban-Ciecko and Barth, 2016), which again is consistent with our model as SST interneurons need to constantly match the top-down input received by neighbouring pyramidal cells. We would predict that these levels of spontaneous firing rates in SST should match the level of feedback received by the pyramidal cells targeted by a particular SST interneuron. In addition, our model predicts the need for a weak top-down input onto SST interneurons. Again, this is in line with recent top-down connectivity studies suggesting that SST can indeed provide such a precise cancellation of the top-down excitatory inputs (Zhang et al., 2014; Leinweber et al., 2017).

(3) *Hierarchy of prediction errors* A further implication of our multi-area learning model is that a high-level prediction error occurring at some higher cortical area would imply also lower-level prediction errors to co-occur at earlier areas. For instance, a categorization error occurring when a visual object is misclassified, would be signalled backwards through our interneuron circuits to lower areas where the individual visual features of the objects are represented. Recent experimental observations in the macaque face-processing hierarchy support this view (Schwiedrzik and Freiwald, 2017). We predict that higher-area activity modulates lower-area activity with the purpose to shape synaptic plasticity at these lower areas.

Here we have focused on the role of SST cells as a feedback-specific interneuron. There are many more interneuron types that we do not consider in our framework. One such type are the PV (parvalbumin-positive) cells, which have been postulated to mediate a somatic excitation-inhibition balance (Vogels et al., 2011; Froemke, 2015) and competition (Masquelier and Thorpe, 2007; Nessler et al., 2013). These functions could in principle be combined with the framework intro-

duced here, or as we suggest below, PV interneuron may be involved in representing yet another type of prediction errors different from the classification errors considered so far. VIP (vasoactive intestinal peptide) interneurons that are believed to be engaged in cortical disinhibition (Letzkus et al., 2015) are assumed in our framework to switch between the discriminative mode and the local attention mode in which lower area activity is generated out of higher area activity (see Fig. 5).

We have focused on an interpretation of our predictive microcircuits as learning across brain areas, but they may also be interpreted as learning across different groups of pyramidal cells within the same brain area.

**Comparison to previous approaches**   It has been suggested that error backpropagation could be approximated by an algorithm that requires alternating between two learning phases, known as contrastive Hebbian learning (Ackley et al., 1985). This link between the two algorithms was first established for an unsupervised learning task (Hinton and McClelland, 1988) and later analyzed (Xie and Seung, 2003) and generalized to a broader class of models (O'Reilly, 1996; Scellier and Bengio, 2017). The two phases needed for contrastive Hebbian learning are: (i) for each input pattern, the network first has to settle down being solely driven by inputs; then, (ii) the process is repeated while additionally driving outputs towards a desired target state. Learning requires subtracting activity patterns recorded on each phase — and therefore requires storing activity in time — or changing plasticity rules across the network on a coordinated, phase-dependent manner, which appears to be biologically implausible.

Two-phase learning recently reappeared in a study which, like ours, uses compartmental neurons (Guerguiev et al., 2017). In this more recent work, the difference between the activity of the apical dendrite in the presence and the absence of the teaching input represents the error that induces plasticity at the forward synapses. This error is used directly for learning the bottom-up synapses without influencing the somatic activity of the pyramidal cell. In contrast, we postulate that the apical dendrite has an explicit error representation at every moment in time by simultaneously integrating top-down excitation and lateral inhibition. As a consequence, we do not need to postulate separate temporal phases, and our network operates continuously in time while plasticity at all synapses is always turned on.

The solution proposed here to avoid two-phase learning relies on a plastic microcircuit that provides functional lateral inhibition. All the involved plasticity rules are error-correcting in spirit and can be understood as learning to match a certain target voltage. For the synapses from the interneu-

rons to the apical dendrite of the pyramidal neurons, the postsynaptic target is the resting potential, and hence the (functionally) inhibitory plasticity rule can be seen as achieving a dendritic balance, similarly to the homeostatic balance of inhibitory synaptic plasticity as previously suggested (Vogels et al., 2011; Luz and Shamir, 2012). Yet, in our model, inhibitory plasticity plays a central role in multi-area, deep error coding, which goes beyond the standard view of inhibitory plasticity as a homeostatic stabilizing force (Keck et al., 2017).

Error minimization is an integral part of brain function according to predictive coding theories (Rao and Ballard, 1999; Friston, 2005), and backprop can be mapped onto a predictive coding network architecture (Whittington and Bogacz, 2017). From a formal point of view this approach is encompassed by the framework introduced by LeCun (1988). A possible network implementation is suggested by Whittington and Bogacz (2017) that requires intricate circuitry with appropriately tuned error-representing neurons. According to that model, the only plastic synapses are those that connect prediction and error neurons.

We built upon the previously made observation that top-down and bottom-up weights need not be in perfect agreement to enable multi-area error-driven learning (Lee et al., 2015; Lillicrap et al., 2016). Consistent with these findings, the strict weight symmetry arising in the classical error backpropagation algorithm is not required in our case either for a successful learning in hidden area neurons.

We have also shown that top-down synapses can be learned using the same dendritic predictive learning rule used at the remaining connections. In our model, the top-down connections have a dual role: they are involved in the apical error representation and, they learn to match the somatic firing driven by the bottom-up input (Urbanczik and Senn, 2014). The simultaneous learning of the bottom-up and top-down pathways leads to the formation of a generative network that can denoise sensory input or generate dream-like inputs (Fig. 5).

Finally, the framework introduced here could also be adapted to other types of error-based learning, such as in generative models that instead of learning to discriminate sensory inputs, learn to generate following sensory input statistics. Error propagation in these forms of generative models, which arise from an inaccurate prediction of sensory inputs, may rely on different dendritic compartments and interneurons, such as the previously mentioned PV inhibitory cells (Petreanu et al., 2009).

## Acknowledgements

## Methods

**Neuron and network model.** The somatic membrane potentials of pyramidal neurons and interneurons evolve in time according to

$$\frac{d}{dt}\mathbf{u}_k^P(t) = -g_{\text{lk}}\,\mathbf{u}_k^P(t) + g_B\left(\mathbf{v}_{B,k}^P(t) - \mathbf{u}_k^P(t)\right) + g_A\left(\mathbf{v}_{A,k}^P(t) - \mathbf{u}_k^P(t)\right) + \sigma\,\xi(t) \tag{1}$$

$$\frac{d}{dt}\mathbf{u}_k^I(t) = -g_{\text{lk}}\,\mathbf{u}_k^I(t) + g_D\left(\mathbf{v}_k^I(t) - \mathbf{u}_k^I(t)\right) + \mathbf{i}_k^I(t) + \sigma\,\xi(t), \tag{2}$$

with one such pair of dynamical equations for every hidden area $0 < k < N$; input area neurons are indexed by $k = 0$.

Eqs. 1 and 2 describe standard conductance-based voltage integration dynamics, having set membrane capacitance to unity and resting potential to zero for clarity purposes. Background activity is modelled as a Gaussian white noise input, $\xi$ in the equations above. To keep the exposition brief we use matrix notation, and denote by $\mathbf{u}_k^P$ and $\mathbf{u}_k^I$ the vectors of pyramidal and interneuron somatic voltages, respectively. Both matrices and vectors, assumed column vectors by default, are typed in boldface here and throughout.

As described in the main text, pyramidal hidden neurons are taken as three-compartment neurons to explicitly incorporate basal and apical dendritic integration zones into the model, inspired by the design of L2/3 pyramidal cells. The two dendritic compartments are coupled to the soma with effective transfer conductances $g_B$ and $g_A$, respectively. Compartmental potentials are given in instantaneous form by

$$\mathbf{v}_{B,k}^P(t) = \mathbf{W}_{k,k-1}^{PP}\,\phi(\mathbf{u}_{k-1}^P(t)) \tag{3}$$

$$\mathbf{v}_{A,k}^P(t) = \mathbf{W}_{k,k+1}^{PP}\,\phi(\mathbf{u}_{k+1}^P(t)) + \mathbf{W}_{k,k}^{PI}\,\phi(\mathbf{u}_k^I(t)), \tag{4}$$

where $\phi(\mathbf{u})$ is the neuronal transfer function, which acts componentwise on $\mathbf{u}$.

Although the design can be extended to more complex morphologies, in the framework of dendritic predictive plasticity two compartments suffice to compare desired target with actual prediction. Hence, aiming for simplicity, we reduce pyramidal output neurons to two-compartment cells, essentially following Urbanczik and Senn (2014); the apical compartment is absent ($g_A = 0$ in Eq. 1) and basal voltages are as defined in Eq. 3. Synapses proximal to the somata of output neurons provide direct external teaching input, incorporated as an additional source of current $\mathbf{i}_N^P$. For any given such neuron, excitatory and inhibitory conductance-based input generates a somatic current

$i_N^P(t) = g_{\text{exc},N}^P(t)\left(E_{\text{exc}} - u_N^P(t)\right) + g_{\text{inh},N}^P(t)\left(E_{\text{inh}} - u_N^P(t)\right)$, where $E_{\text{exc}}$ and $E_{\text{inh}}$ are excitatory and inhibitory synaptic reversal potentials, respectively. The point at which no current flows, $i_N^P = 0$, defines the target teaching voltage $u_N^{\text{trgt}}$ towards which the neuron is nudged.

Interneurons are similarly modelled as two-compartment cells, cf. Eq. 2. Lateral dendritic projections from neighboring pyramidal neurons provide the main source of input

$$\mathbf{v}_k^I(t) = \mathbf{W}_{k,k}^{IP}\, \phi(\mathbf{u}_k^P(t)), \tag{5}$$

whereas cross-area, top-down synapses define the teaching current $\mathbf{i}_k^I$. Specifically, an interneuron at area $k$ receives private somatic teaching excitatory and inhibitory input from a pyramidal neuron at area $k$+1 balanced according to $g_{\text{exc},k}^I = g_{\text{som}}\frac{u_{k+1}^P - E_{\text{inh}}}{E_{\text{exc}} - E_{\text{inh}}}$, $g_{\text{inh},k}^I = -g_{\text{som}}\frac{u_{k+1}^P - E_{\text{exc}}}{E_{\text{exc}} - E_{\text{inh}}}$, where $g_{\text{som}}$ is some constant scale factor denoting overall nudging strength; with this setting, the interneuron is nudged to follow the corresponding next area pyramidal neuron.

**Synaptic plasticity.** Our model synaptic weight update functions belong to the class of dendritic predictive plasticity rules (Urbanczik and Senn, 2014; Spicher et al., in preparation) that can be expressed in general form as

$$\frac{d}{dt}w = \eta\, h(v)\left(\phi(u) - \phi(v)\right) r, \tag{6}$$

where $w$ is an individual synaptic weight, $\eta$ is a learning rate, $u$ and $v$ denote distinct compartmental potentials, $\phi$ is a rate function, third factor $h$ is a function of potential $v$, and $r$ is the presynaptic input. Eq. 6 was originally derived in the light of reducing the prediction error of somatic spiking, when $u$ represents the somatic potential and $v$ is a function of the postsynaptic dendritic potential.

Concretely, the plasticity rules for the various connection types present in the network are:

$$\frac{d}{dt}\mathbf{W}_{k,k-1}^{PP} = \eta_{k,k-1}^{PP}\left(\phi(\mathbf{u}_k^P) - \phi(\hat{\mathbf{v}}_{B,k}^P)\right)\left(\mathbf{r}_{k-1}^P\right)^T, \tag{7}$$

$$\frac{d}{dt}\mathbf{W}_{k,k}^{IP} = \eta_{k,k}^{IP}\left(\phi(\mathbf{u}_k^I) - \phi(\hat{\mathbf{v}}_k^I)\right)\left(\mathbf{r}_k^P\right)^T, \tag{8}$$

$$\frac{d}{dt}\mathbf{W}_{k,k}^{PI} = \eta_{k,k}^{PI}\left(\mathbf{v}_{\text{rest}} - \mathbf{v}_{A,k}^P\right)\left(\mathbf{r}_k^I\right)^T, \tag{9}$$

where $(\cdot)^T$ denotes vector transpose and $\mathbf{r}_k \equiv \phi(\mathbf{u}_k)$ the area $k$ firing rates. So the strengths of plastic synapses evolve according to the correlation of dendritic prediction error and presynaptic rate and can undergo both potentiation or depression depending on the sign of the first factor.

For basal synapses, such prediction error factor amounts to a difference between postsynaptic

rate and a local dendritic estimate which depends on the branch potential. In Eqs. 7 and 8, dendritic predictions $\hat{\mathbf{v}}_{B,k}^P = \frac{g_B}{g_{lk}+g_B+g_A} \mathbf{v}_{B,k}^P$ and $\hat{\mathbf{v}}_k^I = \frac{g_D}{g_{lk}+g_D} \mathbf{v}_k^I$ take into account dendritic attenuation factors. Meanwhile, plasticity rule (9) of lateral interneuron-to-pyramidal synapses aims to silence (i.e., set to resting potential $\mathbf{v}_{rest} = \mathbf{0}$, here and throughout null for simplicity) the apical compartment; this introduces an attractive state for learning where the contribution from interneurons balances top-down dendritic input. The learning rule of apical-targetting synapses can be thought of as a dendritic variant of the homeostatic inhibitory plasticity proposed by Vogels et al. (2011).

In the experiments where the top-down connections are plastic (cf. Fig. 5), the weights evolve according to

$$\frac{d}{dt}\mathbf{W}_{k,k+1}^{PP} = \eta_{k,k+1}^{PP} \left(\phi(\mathbf{u}_k^P) - \phi(\hat{\mathbf{v}}_{TD,k}^P)\right) \left(\mathbf{r}_{k+1}^P\right)^T, \tag{10}$$

with $\hat{\mathbf{v}}_{TD,k}^P = \mathbf{W}_{k,k+1} \mathbf{r}_{k+1}^P$. An implementation of this rule requires a subdivision of the apical compartment into a distal part receiving the top-down input (with voltage $\hat{\mathbf{v}}_{TD,k}^P$) and a more proximal part receiving the lateral input from the interneurons (with voltage $\mathbf{v}_{A,k}^P$).

**Nonlinear function approximation task.** In Fig. 3, a pyramidal neuron network learns to approximate a random nonlinear function implemented by a held-aside feedforward network with the same (30-20-10) dimensions; this ensures that the target function is realizable. One teaching example consists of a randomly drawn input pattern $\mathbf{r}_0^P$ assigned to corresponding target $\mathbf{r}_2^{trgt} = \phi(\mathbf{W}_{2,1}^{trgt} \phi(\mathbf{W}_{1,0}^{trgt} \mathbf{r}_0^P))$. Teacher network weights and input pattern entries are sampled from a uniform distribution $U(-1, 1)$. We choose a soft rectifying nonlinearity as the neuronal transfer function, $\phi(u) = \log(1 + \exp(u))$.

The pyramidal neuron network is initialized to a self-predicting state where $\mathbf{W}_{1,1}^{IP} = \mathbf{W}_{2,1}^{PP}$ and $\mathbf{W}_{1,1}^{PI} = -\mathbf{W}_{1,2}^{PP}$. Top down weight matrix $\mathbf{W}_{1,2}^{PP}$ is fixed and set at random with entries drawn from a uniform distribution. Output area teaching currents $\mathbf{i}_2^P$ are set so as to nudge $\mathbf{u}_2^P$ towards the teacher-generated $\mathbf{u}_2^{trgt}$. Reported error curves are exponential moving averages of the sum of squared errors loss $\|r_2^P - r_2^{trgt}\|^2$ computed after every example on unseen input patterns. Plasticity induction terms given by the right-hand sides of Eqs. 7-9 are low-pass filtered with time constant $\tau_w$ before being definitely consolidated, to dampen fluctuations; synaptic plasticity is kept on throughout. Plasticity and neuron model parameters are given in the accompanying supplementary material.

**MNIST image classification and reconstruction tasks.** When simulating the larger models

used on the MNIST data set we resort to a discrete-time network dynamics where the compartmental potentials are updated in two steps before applying synaptic changes.

The simplified model dynamics is as follows. For each presented MNIST image, both pyramidal and interneurons are first initialized to their bottom-up prediction state (3), $\mathbf{u}_k = \mathbf{v}_{\text{B},k}$, starting from area 1 upto top area $N$. Output area neurons are then nudged towards their desired target $\mathbf{u}_N^{\text{trgt}}$, yielding updated somatic potentials $\mathbf{u}_N^{\text{P}} = (1 - \lambda_N)\,\mathbf{v}_{\text{B},N} + \lambda_N\,\mathbf{u}_N^{\text{trgt}}$. To obtain the remaining final compartmental potentials, the network is revisited in reverse order, proceeding from area $k = N - 1$ down to $k = 1$. For each $k$, interneurons are first updated to include top-down teaching signals, $\mathbf{u}_k^{\text{I}} = (1 - \lambda_I)\,\mathbf{v}_k^{\text{I}} + \lambda_I\,\mathbf{u}_{k+1}^{\text{P}}$; this yields apical compartment potentials according to (4), after which we update hidden area somatic potentials as a convex combination with mixing factor $\lambda_k$. The convex combination factors introduced above are directly related to neuron model parameters as conductance ratios. Synaptic weights are then updated according to Eqs. 7-10.

Such simplified dynamics approximates the full recurrent network relaxation in the deterministic setting $\sigma \to 0$, with the approximation improving as the top-down dendritic coupling is decreased, $g_{\text{A}} \to 0$.

We train the models on the standard MNIST handwritten image database, further splitting the training set into 55000 training and 5000 validation examples. The reported test error curves are computed on the 10000 held-aside test images. The four-area network shown in Fig. 4 is initialized in a self-predicting state with appropriately scaled initial weight matrices. To speed-up training we use a mini-batch strategy on every learning rule, whereby weight changes are averaged across 10 images before being actually consolidated. We take the neuronal transfer function $\phi$ to be a logistic function, $\phi(u) = 1/(1 + \exp(-u))$ and include a learnable threshold on each neuron, modelled as an additional input fixed at unity with plastic weight. Desired target class vectors are 1-hot coded, with $r_N^{\text{trgt}} \in \{0.1, 0.8\}$. During testing, the output is determined by picking the class label corresponding to the neuron with highest firing rate. Model parameters are given in full in the supplementary material.

To generate digit prototypes as shown in Fig. 5C, the network is ran feedforward in a top-to-bottom fashion: a pass of pyramidal neuron activations is performed, while disabling the feedforward stream as well as the interneuron negative lateral contributions. For this reason, this mode of recall is referred to in the main text as the disinhibited mode. The output area is initialized to the 1-hot-coded pattern corresponding to the desired digit class.

The denoised images shown in Fig. 5D are the top-down predictions $\hat{\mathbf{r}}_0 = \phi(\hat{\mathbf{v}}_{\text{TD},0}^{\text{P}})$ obtained after

presenting randomly selected digit examples from the test set, corrupted with additive Gaussian noise of standard deviation $\sigma = 0.3$. The network states are determined by the two-step procedure described above. Recurrent effects are therefore ignored, as a single backward step is performed.

**Computer code.** For the first series of experiments (Figs. 1-3) we wrote custom Mathematica (Wolfram Research, Inc.) code. The larger MNIST networks (Figs. 4 and 5) were simulated in Python using the TensorFlow framework.

# References

Ackley DH, Hinton GE, Sejnowski TJ (1985) A learning algorithm for Boltzmann machines. *Cognitive Science* 9:147–169.

Attinger A, Wang B, Keller GB (2017) Visuomotor coupling shapes the functional development of mouse visual cortex. *Cell* 169:1291–1302.e14.

Bengio Y (2014) How auto-encoders could provide credit assignment in deep networks via target propagation. *arXiv:1407.7906*.

Berger TK, Perin R, Silberberg G, Markram H (2009) Frequency-dependent disynaptic inhibition in the pyramidal network: a ubiquitous pathway in the developing rat neocortex. *The Journal of physiology* 587:5411–5425.

Berger TK, Silberberg G, Perin R, Markram H (2010) Brief bursts self-inhibit and correlate the pyramidal network. *PLOS Biology* 8:e1000473.

Bono J, Clopath C (2017) Modeling somatic and dendritic spike mediated plasticity at the single neuron and network level. *Nature Communications* 8:706.

Bottou L (1998) Online algorithms and stochastic approximations. In Saad D, editor, *Online Learning and Neural Networks*. Cambridge University Press, Cambridge, UK.

Cauller LJ, Kulics AT (1991) The neural basis of the behaviorally relevant N1 component of the somatosensory-evoked potential in SI cortex of awake monkeys: evidence that backward cortical projections signal conscious touch sensation. *Experimental Brain Research* 84:607–619.

Clopath C, Büsing L, Vasilaki E, Gerstner W (2010) Connectivity reflects coding: a model of voltage-based stdp with homeostasis. *Nature Neuroscience* 13:344–352.

Costa RP, Assael YM, Shillingford B, de Freitas N, Vogels TP (2017) Cortical microcircuits as gated-recurrent neural networks In *Advances in Neural Information Processing Systems*, pp. 271–282.

Crick F (1989) The recent excitement about neural networks. *Nature* 337:129–132.

Dorrn AL, Yuan K, Barker AJ, Schreiner CE, Froemke RC (2010) Developmental sensory experience balances cortical excitation and inhibition. *Nature* 465:932–936.

Friedrich J, Urbanczik R, Senn W (2011) Spatio-temporal credit assignment in neuronal population learning. *PLOS Computational Biology* 7:e1002092.

Friston K (2005) A theory of cortical responses. *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 360:815–836.

Froemke RC (2015) Plasticity of cortical excitatory-inhibitory balance. *Annual Review of Neuroscience* 38:195–219.

Grossberg S (1987) Competitive learning: From interactive activation to adaptive resonance. *Cognitive Science* 11:23–63.

Guerguiev J, Lillicrap TP, Richards BA (2017) Towards deep learning with segregated dendrites. *eLife* 6:e22901.

Hinton GE, McClelland JL (1988) Learning representations by recirculation. In Anderson DZ, editor, *Neural Information Processing Systems*, pp. 358–366. American Institute of Physics.

Huber D, Gutnisky DA, Peron S, O'Connor DH, Wiegert JS, Tian L, Oertner TG, Looger LL, Svoboda K (2012) Multiple dynamic representations in the motor cortex during sensorimotor learning. *Nature* 484:473–478.

Keck T, Toyoizumi T, Chen L, Doiron B, Feldman DE, Fox K, Gerstner W, Haydon PG, Hübener M, Lee HK, Lisman JE, Rose T, Sengpiel F, Stellwagen D, Stryker MP, Turrigiano GG, van Rossum MC (2017) Integrating Hebbian and homeostatic plasticity: the current state of the field and future research directions. *Philosophical Transactions of the Royal Society of London B: Biological Sciences* 372.

Khaligh-Razavi SM, Kriegeskorte N (2014) Deep supervised, but not unsupervised, models may explain it cortical representation. *PLOS Computational Biology* 10:1–29.

Larkum M (2013) A cellular mechanism for cortical associations: an organizing principle for the cerebral cortex. *Trends in Neurosciences* 36:141–151.

LeCun Y (1988) A theoretical framework for back-propagation. In Touretzky D, Hinton G, Sejnowski T, editors, *Proceedings of the 1988 Connectionist Models Summer School*, pp. 21–28. Morgan Kaufmann, Pittsburg, PA.

LeCun Y, Bengio Y, Hinton G (2015) Deep learning. *Nature* 521:436–444.

Lee DH, Zhang S, Fischer A, Bengio Y (2015) Difference target propagation. In *Machine Learning and Knowledge Discovery in Databases*, pp. 498–515. Springer.

Leinweber M, Ward DR, Sobczak JM, Attinger A, Keller GB (2017) A Sensorimotor Circuit in Mouse Cortex for Visual Flow Predictions. *Neuron* 95:1420–1432.e5.

Letzkus JJ, Wolff SBE, Lüthi A (2015) Disinhibition, a Circuit Mechanism for Associative Learning and Memory. *Neuron* 88:264–276.

Lillicrap TP, Cownden D, Tweed DB, Akerman CJ (2016) Random synaptic feedback weights support error backpropagation for deep learning. *Nature Communications* 7:13276.

Luz Y, Shamir M (2012) Balancing feed-forward excitation and inhibition via Hebbian inhibitory synaptic plasticity. *PLOS Computational Biology* 8:e1002334.

Makino H, Komiyama T (2015) Learning enhances the relative impact of top-down processing in the visual cortex. *Nature Neuroscience* 18:1116–1122.

Manita S, Suzuki T, Homma C, Matsumoto T, Odagawa M, Yamada K, Ota K, Matsubara C, Inutsuka A, Sato M et al. (2015) A top-down cortical circuit for accurate sensory perception. *Neuron* 86:1304–1316.

Marblestone AH, Wayne G, Kording KP (2016) Toward an integration of deep learning and neuroscience. *Frontiers in Computational Neuroscience* 10:94.

Markram H, Toledo-Rodriguez M, Wang Y, Gupta A, Silberberg G, Wu C (2004) Interneurons of the neocortical inhibitory system. *Nature Reviews Neuroscience* 5:793–807.

Masquelier T, Thorpe S (2007) Unsupervised learning of visual features through spike timing dependent plasticity. *PLOS Computational Biology* 3.

Murayama M, Pérez-Garci E, Nevian T, Bock T, Senn W, Larkum ME (2009) Dendritic encoding of sensory stimuli controlled by deep cortical interneurons. *Nature* 457:1137–1141.

Nessler B, Pfeiffer M, Buesing L, Maass W (2013) Bayesian computation emerges in generic cortical microcircuits through spike-timing-dependent plasticity. *PLOS Computational Biology* 9:e1003037.

O'Reilly RC (1996) Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural Computation* 8:895–938.

Petreanu L, Gutnisky DA, Huber D, Xu Nl, O'Connor DH, Tian L, Looger L, Svoboda K (2012) Activity in motor-sensory projections reveals distributed coding in somatosensation. *Nature* 489:299–303.

Petreanu L, Mao T, Sternson SM, Svoboda K (2009) The subcellular organization of neocortical excitatory connections. *Nature* 457:1142–1145.

Pi HJ, Hangya B, Kvitsiani D, Sanders JI, Huang ZJ, Kepecs A (2013) Cortical interneurons that specialize in disinhibitory control. *Nature* 503:521–524.

Poort J, Khan AG, Pachitariu M, Nemri A, Orsolic I, Krupic J, Bauza M, Sahani M, Keller GB, Mrsic-Flogel TD, Hofer SB (2015) Learning enhances sensory and multiple non-sensory representations in primary visual cortex. *Neuron* 86:1478–1490.

Rao RP, Ballard DH (1999) Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects. *Nature Neuroscience* 2:79–87.

Roelfsema PR, van Ooyen A (2005) Attention-gated reinforcement learning of internal representations for classification. *Neural Computation* 17:2176–2214.

Rumelhart DE, Hinton GE, Williams RJ (1986) Learning representations by back-propagating errors. *Nature* 323:533–536.

Scellier B, Bengio Y (2017) Equilibrium propagation: Bridging the gap between energy-based models and backpropagation. *Frontiers in Computational Neuroscience* 11:24.

Schwiedrzik CM, Freiwald WA (2017) High-level prediction signals in a low-level area of the macaque face-processing hierarchy. *Neuron* 96:89–97.e4.

Silberberg G, Markram H (2007) Disynaptic inhibition between neocortical pyramidal cells mediated by Martinotti cells. *Neuron* 53:735–746.

Spicher D, Clopath C, Senn W (in preparation) Predictive plasticity in dendrites: from a computational principle to experimental data.

Spruston N (2008) Pyramidal neurons: dendritic structure and synaptic integration. *Nature Reviews Neuroscience* 9:206–221.

Sutton RS, Barto AG (1998) *Reinforcement learning: An introduction*, Vol. 1 MIT Press, Cambridge, Mass.

Takahashi N, Oertner TG, Hegemann P, Larkum ME (2016) Active cortical dendrites modulate perception. *Science* 354:1587–1590.

Urban-Ciecko J, Barth AL (2016) Somatostatin-expressing neurons in cortical networks. *Nature Reviews Neuroscience* 17:401–409.

Urbanczik R, Senn W (2014) Learning by the dendritic prediction of somatic spiking. *Neuron* 81:521–528.

Vogels TP, Sprekeler H, Zenke F, Clopath C, Gerstner W (2011) Inhibitory plasticity balances excitation and inhibition in sensory pathways and memory networks. *Science* 334:1569–1573.

Whittington JCR, Bogacz R (2017) An approximation of the error backpropagation algorithm in a predictive coding network with local Hebbian synaptic plasticity. *Neural Computation* 29:1229–1262.

Xie X, Seung HS (2003) Equivalence of backpropagation and contrastive Hebbian learning in a layered network. *Neural Computation* 15:441–454.

Yamins DL, DiCarlo JJ (2016) Using goal-driven deep learning models to understand sensory cortex. *Nature Neuroscience* 19:356–365.

Yamins DL, Hong H, Cadieu CF, Solomon EA, Seibert D, DiCarlo JJ (2014) Performance-optimized hierarchical models predict neural responses in higher visual cortex. *Proceedings of the National Academy of Sciences* 111:8619–8624.

Zhang S, Xu M, Kamigaki T, Do JPH, Chang WC, Jenvay S, Miyamichi K, Luo L, Dan Y (2014) Long-range and local circuits for top-down modulation of visual cortex processing. *Science* 345:660–665.

Zmarz P, Keller GB (2016) Mismatch receptive fields in mouse visual cortex. *Neuron* 92:766–772.

# Supplementary information

## Supplementary data

Below we detail the model parameters used to generate the figures presented in the main text.

**Fig. 1 details**. The parameters for the compartmental model neuron were: $g_A = 0.8$, $g_B = g_D = 1.0$, $g_{lk} = 0.1$. Interneuron somatic teaching conductances were balanced to yield overall nudging strength $g_{som} = 0.8$. Initial weight matrix entries were independently drawn from a uniform distribution $U(-1, 1)$. We chose background activity levels of $\sigma = 0.1$. The learning rates were set as $\eta_{1,1}^{IP} = 0.0002375$ and $\eta_{1,1}^{PI} = 0.0005$.

Input patterns were smoothly transitioned by low-pass filtering $\mathbf{u}_0^P$ with time constant $\tau_0 = 3$. A transition between patterns was triggered every 100 ms. Weight changes were low pass filtered with time constant $\tau_w = 30$. The dynamical equations were solved using Euler's method with a time step of 0.1, which resulted in 1000 integration time steps per pattern.

**Fig. 2 details**. We used learning rates $\eta_{1,0}^{PP} = \eta_{1,1}^{IP} = 0.0011875$ and $\eta_{2,1}^{PP} = 0.0005$. Remaining parameters as used for Fig. 1.

**Fig. 3 details**. Initial forward weights $\mathbf{W}_{1,0}^{PP}$ and $\mathbf{W}_{2,1}^{PP}$ were scaled down by a factor of 0.1. Background noise level was raised to $\sigma = 0.3$. The learning rates were $\eta_{1,1}^{IP} = 0.00002375$, $\eta_{1,0}^{PP} = 0.00011875$, $\eta_{2,1}^{PP} = 0.00001$. Weight matrices $\mathbf{W}_{1,2}^{PP}$ and $\mathbf{W}_{1,1}^{PI}$ were kept fixed, so the model relied on a feedback alignment mechanism to learn. Remaining parameters as used for Fig. 1.

**Fig. 4 details**. We chose mixing factors $\lambda_3 = \lambda_I = 0.1$ and $\lambda_1 = \lambda_2 = 0.3$. Forward learning rates were $\eta_{3,2}^{PP} = 0.001/\lambda_3$, $\eta_{2,1}^{PP} = \eta_{3,2}^{PP}/\lambda_2$, $\eta_{1,0}^{PP} = \eta_{2,1}^{PP}/\lambda_1$. Lateral learning rates were $\eta_{2,2}^{IP} = 2\eta_{3,2}^{PP}$ and $\eta_{1,1}^{IP} = 2\eta_{2,1}^{PP}$. Initial forward weights were drawn at random from a uniform distribution $U(-0.1, 0.1)$, and the remaining weights from $U(-1, 1)$.

**Fig. 5 details**. We took all mixing factors equal $\lambda_2 = \lambda_1 = \lambda_I = 0.1$. Forward learning rates: $\eta_{2,1}^{PP} = 0.02/\lambda_2$, $\eta_{1,0}^{PP} = \eta_{2,1}^{PP}/\lambda_1$. Lateral connections learned with rate $\eta_{1,1}^{IP} = \eta_{1,1}^{PI} = \eta_{2,1}^{PP}$. Top-down connections were initialized from a uniform distribution $U(-0.1, 0.1)$ and adapted with learning rates $\eta_{1,2}^{PP} = 0.0002$ and $\eta_{0,1}^{PP} = 0.0001$.

## Supplementary analysis

In this supplementary note we present a set of mathematical results concerning the network and plasticity model described in the main text.

To proceed analytically we make a number of simplifying assumptions. Unless noted otherwise, we study the network in a deterministic setting and consider the limiting case where lateral microcircuit synaptic weights match the corresponding forward weights:

$$\mathbf{W}^{\text{PI}}_{k,k} = -\mathbf{W}^{\text{PP}}_{k,k+1} \equiv \mathbf{W}^{\text{PI}*}_{k,k} \tag{S1}$$

$$\mathbf{W}^{\text{IP}}_{k,k} = \frac{g_{\text{B}} + g_{\text{lk}}}{g_{\text{B}} + g_{\text{A}} + g_{\text{lk}}} \mathbf{W}^{\text{PP}}_{k+1,k} \equiv \mathbf{W}^{\text{IP}*}_{k,k}, \tag{S2}$$

The particular choice of proportionality factors, which depend on the neuron model parameters, is motivated below. Under the above configuration, the network becomes self-predicting.

To formally relate the encoding and propagation of errors implemented by the inhibitory microcircuit to the backpropagation of errors algorithm from machine learning, we consider the limit where top-down input is weak compared to the bottom-up drive. This limiting case results in error signals that decrease exponentially with area depth, but allows us to proceed analytically.

We further assume that the top-down weights converging to the apical compartments are equal to the corresponding forward weights, $\mathbf{W}^{\text{PP}}_{k,k+1} = \left(\mathbf{W}^{\text{PP}}_{k+1,k}\right)^{T}$. Such weight symmetry is not essential for successful learning in a broad range of problems, as demonstrated in the main simulations and as observed before (Lee et al., 2015; Lillicrap et al., 2016). It is, however, required to frame learning as a gradient descent procedure. Furthermore, in the analyses of the learning rules, we assume that synaptic changes take place at a fixed point of the neuronal dynamics; we therefore consider discrete-time versions of the plasticity rules. This approximates the continuous-time plasticity model as long as changes in the inputs are slow compared to the neuronal dynamics.

For convenience, we will occasionally drop neuron type indices and refer to bottom-up weights $\mathbf{W}_{k+1,k}$ and to top-down weights $\mathbf{W}_{k,k+1}$. Additionally, we assume without loss of generality that the dendritic coupling conductance for interneurons is equal to the basal dendritic coupling of pyramidal neurons, $g_{\text{D}} = g_{\text{B}}$. Finally, whenever it is useful to distinguish whether output area nudging is turned off, we use superscript '$-$'.

**Interneuron activity in the self-predicting state.** Following Urbanczik and Senn (2014), we note that steady state interneuron somatic potentials can be expressed as a convex combination of

basal dendritic and pyramidal neuron potentials that are provided via somatic teaching input:

$$\mathbf{u}_k^{\mathsf{I}} = \frac{g_{\mathsf{B}}}{g_{\mathsf{lk}} + g_{\mathsf{B}} + g_{\mathsf{som}}} \mathbf{v}_k^{\mathsf{I}} + \frac{g_{\mathsf{som}}}{g_{\mathsf{lk}} + g_{\mathsf{B}} + g_{\mathsf{som}}} \mathbf{u}_{k+1}^{\mathsf{P}} = (1 - \lambda)\,\hat{\mathbf{v}}_k^{\mathsf{I}} + \lambda\,\mathbf{u}_{k+1}^{\mathsf{P}}, \tag{S3}$$

with $g_{\mathsf{B}}$ and $g_{\mathsf{lk}}$ the effective dendritic transfer and leak conductances, respectively, and $g_{\mathsf{som}}$ the total excitatory and inhibitory teaching conductance. In the equation above, $\hat{\mathbf{v}}_k^{\mathsf{I}} = \frac{g_{\mathsf{B}}}{g_{\mathsf{lk}}+g_{\mathsf{B}}}\mathbf{v}_k^{\mathsf{I}}$ is the interneuron dendritic prediction (cf. Eq. 8), and $\lambda \equiv \frac{g_{\mathsf{som}}}{g_{\mathsf{lk}}+g_{\mathsf{B}}+g_{\mathsf{som}}} \in [0, 1[$ is a mixing factor which controls the nudging strength for the interneurons. In other words, the current prediction $\hat{\mathbf{v}}_k^{\mathsf{I}}$ and the teaching signal are averaged with coefficients determined by normalized conductances. We will later consider the weak nudging limit of $\lambda \to 0$.

The relation $\hat{\mathbf{v}}_k^{\mathsf{I}} = \hat{\mathbf{v}}_{\mathsf{B},k+1}^{\mathsf{P}}$ holds when pyramidal-to-interneuron synaptic weights are equal to pyramidal-pyramidal forward weights, up to a scale factor: $\mathbf{W}_{k,k}^{\mathsf{IP}} = \frac{g_{\mathsf{lk}}+g_{\mathsf{B}}}{g_{\mathsf{lk}}+g_{\mathsf{B}}+g_{\mathsf{A}}}\mathbf{W}_{k+1,k}^{\mathsf{PP}}$, which simplifies to $\mathbf{W}_{N-1,N-1}^{\mathsf{IP}} = \mathbf{W}_{N,N-1}^{\mathsf{PP}}$ for the last area where $g_{\mathsf{A}} = 0$ (to reduce clutter, we use the slightly abusive notation whereby $g_{\mathsf{A}}$ should be understood to be zero when referring to output area neurons). This is the reason for the particular choice of ideal pyramidal-to-interneuron weights presented in the preamble. The network is then internally consistent, in the sense that the interneurons predict the model's own predictions, held by pyramidal neurons.

**Bottom-up predictions in the absence of external nudging.** We first study the situation where the input pattern $\mathbf{r}_0$ is stationary and the output area teaching input is disabled, $\mathbf{i}_N^{\mathsf{P}} = 0$. We show that the fixed point of the network dynamics is a state where somatic voltages are equal to basal voltages, up to a dendritic attenuation factor. In other words, the network effectively behaves as if it were feedforward, in the sense that it computes the same function as the corresponding network with equal bottom-up but no top-down or lateral connections.

Specifically, in the absence of external nudging (indicated by the $-$ in the superscript), the somatic voltages of pyramidal and interneuron are given by the bottom-up dendritic predictions,

$$\mathbf{u}_k^{\mathsf{P},-} = \hat{\mathbf{v}}_{\mathsf{B},k}^{\mathsf{P},-} \equiv \frac{g_{\mathsf{B}}}{g_{\mathsf{lk}} + g_{\mathsf{B}} + g_{\mathsf{A}}}\,\mathbf{W}_{k,k-1}^{\mathsf{PP}}\,\phi(\hat{\mathbf{v}}_{\mathsf{B},k-1}^{\mathsf{P},-}) \tag{S4}$$

$$\mathbf{u}_k^{\mathsf{I},-} = \hat{\mathbf{v}}_k^{\mathsf{I},-} \equiv \frac{g_{\mathsf{B}}}{g_{\mathsf{lk}} + g_{\mathsf{B}}}\,\mathbf{W}_{k,k}^{\mathsf{IP}}\,\phi(\hat{\mathbf{v}}_{\mathsf{B},k}^{\mathsf{P},-}). \tag{S5}$$

To show that Eq. S4 describes the state of the network, we start at the output area and set Eq. 1 to zero. Because nudging is turned off, we observe that $\mathbf{u}_N^{\mathsf{P}}$ is equal to $\hat{\mathbf{v}}_{\mathsf{B},N}^{\mathsf{P},-}$ if area $N - 1$ also satisfies $\mathbf{u}_{N-1}^{\mathsf{P}} = \hat{\mathbf{v}}_{\mathsf{B},N-1}^{\mathsf{P},-}$. The same recursively applies to the hidden area below when its

apical voltage vanishes, $\mathbf{v}^{\mathrm{P}}_{\mathrm{A},N-1} = 0$. Now we note that at the fixed point the interneuron cancels the corresponding pyramidal neuron, due to the assumption that the network is in a self-predicting state, which yields $\mathbf{u}^{\mathrm{I}}_{N-1} = \mathbf{u}^{\mathrm{P}}_{N}$. Together with the fact that $\mathbf{W}^{\mathrm{PI}}_{N-1,N-1} = -\mathbf{W}^{\mathrm{PP}}_{N-1,N}$, we conclude that the interneuron contribution to the apical compartment cancels the top-down pyramidal neuron input, yielding the required condition $\mathbf{v}^{\mathrm{P}}_{\mathrm{A},N-1} = 0$.

The above argument can be iterated down to the input area, which is constant, and we arrive at Eq. S4.

**Zero plasticity induction in the absence of nudging.** In view of Eq. S4, which states that in the absence of external nudging the somatic voltages correspond to the basal predictions, no synaptic changes are induced in basal synapses on the pyramidal and interneurons as defined by the plasticity rules (7) and (8), respectively. Similarly, the apical voltages are equal to rest, $\mathbf{v}^{\mathrm{P},-}_{\mathrm{A},k} = \mathbf{v}_{\mathrm{rest}}$, when the top-down input is fully predicted, and no synaptic plasticity is induced in the inter-to-pyramidal neuron synapses, see (9). When noisy background currents are present, the average prediction error is zero, while momentary fluctuations will still trigger plasticity. Note that the above holds when the dynamics is away from equilibrium, under the additional constraint that the integration time constant of interneurons matches that of pyramidal neurons.

**Recursive prediction error propagation.** Prediction errors arise in the model whenever lateral interneurons cannot fully explain top-down input, leading to a deviation from baseline in apical dendrite activity. Here, we look at the network steady state equations for a stationary input pattern $\mathbf{r}_0$ and derive an iterative relationship which establishes the propagation across the network of prediction mismatches originating downstream. The following compartmental potentials are thus evaluated at a fixed point of the neuronal dynamics.

Under the assumption (S1) of matching interneuron-to-pyramidal top-down weights, apical compartment potentials simplify to

$$\mathbf{v}^{\mathrm{P}}_{\mathrm{A},k} = \mathbf{W}_{k,k+1}\left[\phi(\mathbf{u}^{\mathrm{P}}_{k+1}) - \phi(\mathbf{u}^{\mathrm{I}}_{k})\right] = \mathbf{W}_{k,k+1}\,\mathbf{e}_{k+1}, \tag{S6}$$

where we introduced error vector $\mathbf{e}_{k+1}$ defined as the difference between pyramidal and interneuron firing rates. Such deviation can be intuitively understood as an area-wise interneuron prediction mismatch, being zero when interneurons perfectly explain pyramidal neuron activity. We now evaluate this difference vector at a fixed point to obtain a recurrence relation that links consecutive areas.

31

The steady-state somatic potentials of hidden pyramidal neurons are given by

$$\mathbf{u}_k^P = \frac{g_B}{g_{lk} + g_B + g_A} \mathbf{v}_{B,k}^P + \frac{g_A}{g_{lk} + g_B + g_A} \mathbf{v}_{A,k}^P = \hat{\mathbf{v}}_{B,k}^P + \lambda \mathbf{v}_{A,k}^P$$

$$= \hat{\mathbf{v}}_{B,k}^P + \lambda \mathbf{W}_{k,k+1} \mathbf{e}_{k+1}. \tag{S7}$$

To shorten the following, we assumed that the apical attenuation factor is equal to the interneuron nudging strength $\lambda$. As previously mentioned, we proceed under the assumption of weak feedback, $\lambda$ small. As for the corresponding interneurons, we insert Eq. S7 into Eq. S3 and note that when the network is in a self-predicting state we have $\hat{\mathbf{v}}_{k-1}^I = \hat{\mathbf{v}}_{B,k}^P$, yielding

$$\mathbf{u}_{k-1}^I = (1 - \lambda)\,\hat{\mathbf{v}}_{B,k}^P + \lambda\left(\hat{\mathbf{v}}_{B,k}^P + \lambda\,\mathbf{v}_{A,k}^P\right) = \hat{\mathbf{v}}_{B,k}^P + \lambda^2\,\mathbf{v}_{A,k}^P. \tag{S8}$$

Using the identities (S7) and (S8), we now expand to first order the difference vector $\mathbf{e}_k$ around $\hat{\mathbf{v}}_{B,k}^P$ as follows

$$\mathbf{e}_k = \phi(\mathbf{u}_k^P) - \phi(\mathbf{u}_{k-1}^I) = \lambda\,\mathbf{D}_k\,\mathbf{v}_{A,k}^P + \mathcal{O}\!\left(\lambda^2\,\|\mathbf{v}_{A,k}^P\|\right). \tag{S9}$$

Matrix $\mathbf{D}_k$ is a diagonal matrix with diagonal equal to $\phi'(\hat{\mathbf{v}}_{B,k}^P)$, i.e., whose $i$-th element reads $\frac{d\phi}{dv}(\hat{v}_{B,k,i}^P)$. It contains the derivative of the neuronal transfer function $\phi$ evaluated component-wise at the bottom-up predictions $\hat{\mathbf{v}}_{B,k+1}^P$. Recalling Eq. S6, we obtain a recurrence relation

$$\mathbf{e}_k = \lambda\,\mathbf{D}_k\,\mathbf{W}_{k,k+1}\,\mathbf{e}_{k+1} + \mathcal{O}\!\left(\lambda^2\,\|\mathbf{W}_{k,k+1}\,\mathbf{e}_{k+1}\|\right). \tag{S10}$$

Finally, last area pyramidal neurons provide the initial condition by being directly nudged towards the desired target $\mathbf{u}_N^{\text{trgt}}$. Their membrane potentials can be written as

$$\mathbf{u}_N^P = (1 - \lambda)\,\hat{\mathbf{v}}_{B,N}^P + \lambda\,\mathbf{u}_N^{\text{trgt}}, \tag{S11}$$

and this gives an estimate for the error in the output area of the form

$$\mathbf{e}_N = \lambda\,\mathbf{D}_N\left(\mathbf{u}_N^{\text{trgt}} - \hat{\mathbf{v}}_{B,N}^P\right) + \mathcal{O}\!\left(\lambda^2\,\|\mathbf{u}_N^{\text{trgt}} - \hat{\mathbf{v}}_{B,N}^P\|\right), \tag{S12}$$

where for simplicity we took the same mixing factor $\lambda$ for pyramidal output and interneurons. Then, for an arbitrary area, assuming that the synaptic weights and the remaining fixed parameters do not

scale with $\lambda$, we arrive at

$$\mathbf{e}_k = \lambda^{N-k+1} \left( \prod_{l=k}^{N-1} \mathbf{D}_l \, \mathbf{W}_{l,l+1} \right) \mathbf{D}_N \left( \mathbf{u}_N^{\text{trgt}} - \hat{\mathbf{v}}_{\text{B},N}^{\text{P}} \right) + \mathcal{O}(\lambda^{N-k+2}). \tag{S13}$$

Thus, steady state potentials of apical dendrites (cf. Eq. S6) recursively encode neuron-specific prediction errors that can be traced back to a mismatch at the output area.

**Learning as approximate error backpropagation.** In the previous section we found that neurons implicitly carry and transmit error information across the network. We now show how the proposed synaptic plasticity model, when applied at a steady state of the neuronal dynamics, can be recast as an approximate gradient descent learning procedure.

More specifically, we compare our model against learning through backprop (Rumelhart et al., 1986) or approximations thereof (Lee et al., 2015; Lillicrap et al., 2016) the weights of the feedfoward multi-area network obtained by removing interneurons and top-down connections from the intact network. For this reference model, the activations $\mathbf{u}_k^-$ are by construction equal to the bottom-up predictions obtained in the full model when output nudging is turned off, $\mathbf{u}_k^- \equiv \hat{\mathbf{v}}_{\text{B},k}^{\text{P},-}$, cf. Eq. S4. Thus, optimizing the weights in the feedforward model is equivalent to optimizing the predictions of the full model.

Define the loss function

$$\mathcal{L}\left( \mathbf{u}_N^-, \mathbf{u}_N^{\text{trgt}} \right) = - \sum_{i=1}^{N_N} \int_0^{u_{N,i}^-} \phi \left( (1-\lambda)\,\nu + \lambda\,u_{N,i}^{\text{trgt}} \right) - \phi(\nu)\,d\nu, \tag{S14}$$

where $N_N$ denotes the number of output neurons. $\mathcal{L}$ can be thought of as the multi-area, multi-output unit analogue of the loss function optimized by the single neuron model (Urbanczik and Senn, 2014), where it stems directly from the particular chosen form of the learning rule (7). The nudging strength parameter $\lambda \in [0, 1[$ allows controlling the mixing with the target and can be understood as an additional learning rate parameter. Albeit unusual in form, function $\mathcal{L}$ imposes a cost similar to an ordinary squared error loss. Importantly, it has a minimum when $\mathbf{u}_N^- = \mathbf{u}_N^{\text{trgt}}$ and it is lower bounded. Furthermore, it is differentiable with respect to compartmental voltages (and synaptic weights). It is therefore suitable for gradient descent optimization. As a side remark, $\mathcal{L}$ integrates to a quadratic function when $\phi$ is linear.

Gradient descent proceeds by changing synaptic weights according to

$$\Delta \mathbf{W}_{k,k-1} = -\eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}_{k,k-1}}. \tag{S15}$$

The required partial derivatives can be efficiently computed by the backpropagation of errors algorithm. For the network architecture we study, this yields a learning rule of the form

$$\Delta \mathbf{W}_{k,k-1}^{\mathrm{bp}} = \eta \, \mathbf{e}_k^- \, \phi(\mathbf{u}_{k-1}^-)^T. \tag{S16}$$

The error factor $\mathbf{e}_k^-$ can be expressed recursively as follows:

$$\mathbf{e}_k^- = \begin{cases} \phi\left((1-\lambda)\mathbf{u}_N^- + \lambda \mathbf{u}_N^{\mathrm{trgt}}\right) - \phi(\mathbf{u}_N^-) & \text{if } k = N, \\ \mathbf{D}_k^- \mathbf{W}_{k+1,k}^T \mathbf{e}_{k+1}^- & \text{otherwise,} \end{cases} \tag{S17}$$

ignoring constant factors that depend on conductance ratios, which can be dealt with by redefining learning rates or backward pass weights. As in the previous section, matrix $\mathbf{D}_k^-$ is a diagonal matrix, with diagonal equal to $\phi'(\mathbf{u}_k^-)$.

We first compare the fixed point equations of the original network to the feedforward activations of the reference model. Starting from the bottommost hidden area, using Eqs. S6, S7 and S13, we notice that $\mathbf{u}_1^P = \mathbf{u}_1^- + \lambda \mathbf{v}_{A,1}^P = \mathbf{u}_1^- + \mathcal{O}(\lambda^N)$, as the bottom-up input is the same in both cases. Inserting this into second hidden area steady state potentials and linearizing the neuronal transfer function gives $\mathbf{u}_2^P = \mathbf{u}_2^- + \lambda \mathbf{v}_{A,2}^P + \mathcal{O}(\lambda^N) = \mathbf{u}_2^- + \mathcal{O}(\lambda^{N-1})$. This can be repeated and for an arbitrary area and neuron type we find

$$\mathbf{u}_k^P = \mathbf{u}_k^- + \lambda \mathbf{v}_{A,k}^P + \mathcal{O}(\lambda^{N-k+2}) = \mathbf{u}_k^- + \mathcal{O}(\lambda^{N-k+1}) \tag{S18}$$

$$\mathbf{u}_{k-1}^I = \mathbf{u}_k^- + \mathcal{O}(\lambda^{N-k+2}). \tag{S19}$$

Writing Eq. S18 in the first form emphasizes that the apical contributions dominate $\mathcal{O}(\lambda \mathbf{v}_{A,k}^P) = \mathcal{O}(\lambda^{N-k+1})$ the bottom-up corrections, which are of order $\mathcal{O}(\lambda^{N-k+2})$.

Next, we prove that up to a factor and to first order the apical term in Eq. S18 represents the backpropagated error in the feedforward network, $\mathbf{e}_k^-$. Starting from the topmost hidden area apical potentials, we reevaluate difference vector (S12) using (S18). Linearization of the neuronal transfer

function gives

$$\mathbf{v}_{A,N-1}^{P} = \lambda\, \mathbf{W}_{N-1,N}\, \mathbf{D}_{N}^{-} \left( \mathbf{u}_{N}^{\text{trgt}} - \mathbf{u}_{N}^{-} \right) + \mathcal{O}(\lambda^2). \tag{S20}$$

Inserting the expression above into Eq. S18 and using Eq. S19 the apical compartment potentials at area $N-1$ can then be recomputed. This procedure can be iterated until the input area is reached. In general form, somatic membrane potentials at hidden area $k$ can be expressed as

$$\mathbf{u}_{k}^{P} = \mathbf{u}_{k}^{-} + \lambda\, \mathbf{v}_{A,k}^{P} + \mathcal{O}(\lambda^{N-k+2}) \tag{S21}$$

$$= \mathbf{u}_{k}^{-} + \lambda^{N-k+1}\, \mathbf{W}_{k,k+1} \left( \prod_{l=k+1}^{N-1} \mathbf{D}_{l}^{-}\, \mathbf{W}_{l,l+1} \right) \mathbf{D}_{N}^{-} \left( \mathbf{u}_{N}^{\text{trgt}} - \mathbf{u}_{N}^{-} \right) + \mathcal{O}(\lambda^{N-k+2}). \tag{S22}$$

This equation shows that, to leading order of $\lambda$, hidden neurons mix and propagate forward purely bottom-up predictions with top-down errors that are computed at the output area and spread backwards.

We are now in position to compare model synaptic weight updates to the ones prescribed by backprop. Output area updates are exactly equal by construction, $\Delta \mathbf{W}_{N,N-1} = \Delta \mathbf{W}_{N,N-1}^{\text{bp}}$. For pyramidal-to-pyramidal neuron synapses from hidden area $k-1$ to area $k$, we obtain

$$\Delta \mathbf{W}_{k,k-1} = \eta_{k,k-1} \left[ \phi(\mathbf{u}_{k}^{P}) - \phi(\hat{\mathbf{v}}_{B,k}^{P}) \right] \left( \mathbf{r}_{k-1}^{P} \right)^{T}$$

$$= \eta_{k,k-1} \left[ \phi\left( \mathbf{u}_{k}^{-} + \lambda\, \mathbf{v}_{A,k}^{P} + \mathcal{O}(\lambda^{N-k+2}) \right) - \phi(\mathbf{u}_{k}^{-}) \right] \left( \mathbf{r}_{k-1}^{-} + \mathcal{O}(\lambda^{N-k+2}) \right)^{T}$$

$$= \eta_{k,k-1}\, \lambda^{N-k+1} \left( \prod_{l=k}^{N-1} \mathbf{D}_{l}^{-}\, \mathbf{W}_{l,l+1} \right) \mathbf{D}_{N}^{-} \left( \mathbf{u}_{N}^{\text{trgt}} - \mathbf{u}_{N}^{-} \right) \left( \mathbf{r}_{k-1}^{-} \right)^{T} + \mathcal{O}(\lambda^{N-k+2}), \tag{S23}$$

while backprop learning rule (S16) can be written as

$$\Delta \mathbf{W}_{k,k-1}^{\text{bp}} = \eta\, \lambda \left( \prod_{l=k}^{N-1} \mathbf{D}_{l}^{-}\, \mathbf{W}_{l,l+1} \right) \mathbf{D}_{N}^{-} \left( \mathbf{u}_{N}^{\text{trgt}} - \mathbf{u}_{N}^{-} \right) \left( \mathbf{r}_{k-1}^{-} \right)^{T} + \mathcal{O}(\lambda^{2}), \tag{S24}$$

where we used that, to first order, the output area error factor is $\mathbf{e}_{N}^{-} = \lambda\, \mathbf{D}_{N}^{-} \left( \mathbf{u}_{N}^{\text{trgt}} - \mathbf{u}_{N}^{-} \right) + \mathcal{O}(\lambda^{2})$. Hence, up to a factor of $\lambda^{N-k}$ which can be absorbed in the learning rate $\eta_{k,k-1}$, changes induced by synaptic plasticity are equal to the backprop learning rule (S16) in the limit $\lambda \to 0$, provided that the top-down weights are set to the transpose of the corresponding feedforward weights, $\mathbf{W}_{k,k+1} = \mathbf{W}_{k+1,k}^{T}$. The 'quasi-feedforward' condition $\lambda \to 0$ has also been invoked to relate backprop to two-phase contrastive Hebbian learning in Hopfield networks (Xie and Seung, 2003).

In our simulations, top-down weights are either set at random and kept fixed, in which case

Eq. S23 shows that the plasticity model optimizes the predictions according to an approximation of backprop known as feedback alignment (Lillicrap et al., 2016); or learned so as to minimize an inverse reconstruction loss, in which case the network implements a form of difference target propagation (Lee et al., 2015).

**Interneuron plasticity.** The analyses of the previous sections relied on the assumption that the synaptic weights to and from interneurons were set to their ideal values, cf. Eqs. S1 and S2. We now study the plasticity of the lateral microcircuit synapses and show that, under mild conditions, learning rules (8) and (9) yield the desired synaptic weight matrices.

We first study the learning of pyramidal-to-interneuron synapses $\mathbf{W}_{k,k}^{\text{IP}}$. To quantify the degree to which the weights deviate from their optimal setting, we introduce the convex loss function

$$\mathcal{L}_k^{\text{IP}} = \frac{1}{2} \operatorname{Tr} \left\{ (\mathbf{W}_{k,k}^{\text{IP}*} - \mathbf{W}_{k,k}^{\text{IP}})^T (\mathbf{W}_{k,k}^{\text{IP}*} - \mathbf{W}_{k,k}^{\text{IP}}) \right\}, \tag{S25}$$

where $\operatorname{Tr}(\mathbf{M})$ denotes the trace of matrix $\mathbf{M}$ and $\mathbf{W}_{k,k}^{\text{IP}*} = \frac{g_\text{B}+g_\text{lk}}{g_\text{B}+g_\text{A}+g_\text{lk}} \mathbf{W}_{k+1,k}^{\text{PP}}$, as defined in Eq. S2.

Starting from the pyramidal-to-interneuron synaptic plasticity rule (8), we express the interneuron somatic potential in convex combination form (S3) and then expand to first order around $\hat{\mathbf{v}}_k^\text{I}$,

$$\begin{aligned}
\Delta \mathbf{W}_{k,k}^{\text{IP}} &= \eta_{k,k}^{\text{IP}} \left( \phi(\mathbf{u}_k^\text{I}) - \phi(\hat{\mathbf{v}}_k^\text{I}) \right) (\mathbf{r}_k^\text{P})^T \\
&= \eta_{k,k}^{\text{IP}} \lambda \, \mathbf{D}_k^{\text{IP}} \left( \mathbf{u}_{k+1}^\text{P} - \hat{\mathbf{v}}_k^\text{I} \right) (\mathbf{r}_k^\text{P})^T + \mathcal{O}(\lambda^2) \\
&= \eta_{k,k}^{\text{IP}} \lambda \, \frac{g_\text{B}}{g_\text{lk} + g_\text{B}} \mathbf{D}_k^{\text{IP}} \left( \mathbf{W}_{k,k}^{\text{IP}*} - \mathbf{W}_{k,k}^{\text{IP}} \right) \mathbf{Q}_k + \mathcal{O}(\lambda^2) + \mathcal{O}(\lambda \, \alpha).
\end{aligned} \tag{S26}$$

Matrix $\mathbf{Q}_k = \mathbf{r}_k^\text{P} (\mathbf{r}_k^\text{P})^T$ denotes the outer product, and $\mathbf{D}_k^{\text{IP}}$ is a diagonal matrix with $i$-th diagonal entry equal to $\phi'(\hat{\mathbf{v}}_{k,i}^\text{I})$.

For simplicity, we ignore fluctuations arising from the stochastic sequential presentation of patterns (Bottou, 1998) and look only at the expected synaptic dynamics[1]. We absorb irrelevant scale factors and to avoid a vanishing update we rescale the learning rate $\hat{\eta}_{k,k}^{\text{IP}}$ by $\lambda^{-1}$. Then, taking the limit $\lambda \to 0$, $\alpha \to 0$ as in the previous sections yields

$$\begin{aligned}
\operatorname{E}\left[\Delta \mathbf{W}_{k,k}^{\text{IP}}\right] &= \hat{\eta}_{k,k}^{\text{IP}} \left( \mathbf{W}_{k,k}^{\text{IP}*} - \mathbf{W}_{k,k}^{\text{IP}} \right) \operatorname{E}\left[\mathbf{D}_k^{\text{IP}} \mathbf{Q}_k\right] \\
&= -\hat{\eta}_{k,k}^{\text{IP}} \frac{\partial \mathcal{L}_k^{\text{IP}}}{\partial \mathbf{W}_{k,k}^{\text{IP}}} \operatorname{E}\left[\mathbf{D}_k^{\text{IP}} \mathbf{Q}_k\right].
\end{aligned} \tag{S27}$$

---

[1] This can be understood as a batch learning protocol, where weight changes are accumulated in the limit of many patterns before being effectively consolidated as a synaptic update.

The expectation is taken over the pattern ensemble. In the last equality above, we used the fact that the gradient of $\mathcal{L}_k^{\text{IP}}$ with respect to the lateral weights $\mathbf{W}_{k,k}^{\text{IP}}$ is given by the difference $\mathbf{W}_{k,k}^{\text{IP}*} - \mathbf{W}_{k,k}^{\text{IP}}$.

As long as the expectation on the right-hand side of (S27) is positive definite, the synaptic dynamics is within $90^{\text{o}}$ of the gradient and thus leads to the unique minimum of $\mathcal{L}_k^{\text{IP}}$. This condition is easily met in practice. For linear neurons, it amounts to requiring that the correlation matrix $\text{E}[\mathbf{Q}]$ is positive definite. In other words, the patterns have to span $\mathbb{R}^{N_k}$, with $N_k$ being the number of pyramidal neurons at area $k$. This is fulfilled when uncorrelated background noise currents are present, and it is likely the case for deterministic networks solving nontrivial tasks. For nonlinear neurons with a monotonically increasing transfer function $\phi$, saturation can lead to a matrix numerically close to singular and slow down learning. A weight matrix initialization that sets the neurons operating far from saturation is therefore an appropriate choice.

A mathematical analysis of the coupled system defined by the various plasticity rules acting in concert is rather involved and beyond the scope of this note. However, the learning of apical-targetting interneuron-to-pyramidal synapses can be studied in isolation by invoking a separation of timescales argument. To proceed, we assume that pyramidal-to-interneuron synapses are ideally set, $\mathbf{W}_{k,k}^{\text{IP}} = \mathbf{W}_{k,k}^{\text{IP}*}$. In practice, this translates to a choice of a small effective learning rate for apical-targetting weights $\eta_{k,k}^{\text{PI}}$. Note that, indirectly, this requirement also imposes a constraint on how fast top-down pyramidal-to-pyramidal weights $\mathbf{W}_{k,k+1}^{\text{PP}}$ can evolve. This is the parameter regime explored in the main text simulations with plastic top-down weights, Fig. 5.

We can then proceed in a similar manner to the previous analysis, as we briefly outline below. Recalling from Eq. S1 that $\mathbf{W}_{k,k}^{\text{PI}*} = -\mathbf{W}_{k,k+1}^{\text{PP}}$, we define the loss function

$$\mathcal{L}_k^{\text{PI}} = \frac{1}{2}\,\text{Tr}\left\{(\mathbf{W}_{k,k}^{\text{PI}*} - \mathbf{W}_{k,k}^{\text{PI}})^T(\mathbf{W}_{k,k}^{\text{PI}*} - \mathbf{W}_{k,k}^{\text{PI}})\right\}. \tag{S28}$$

After some manipulation, as $\lambda \to 0$ the expected synaptic change can be written as

$$\text{E}\left[\Delta\mathbf{W}_{k,k}^{\text{PI}}\right] = -\eta_{k,k}^{\text{PI}}\,\frac{\partial\mathcal{L}_k^{\text{PI}}}{\partial\mathbf{W}_{k,k}^{\text{PI}}}\,\text{E}\left[\mathbf{r}_{k+1}^{\text{P}}\,(\mathbf{r}_{k+1}^{\text{P}})^T\right], \tag{S29}$$

which leads us to conclude that the weights converge to the appropriate values, provided that the correlation matrix of area $k+1$ activity patterns is positive definite.