

Layer-skipping connections facilitate training of layered networks using equilibrium propagation.

Jimmy Gammell^{1,*}, Sae Woo Nam¹ and Adam N. McCaughan¹

¹National Institute of Standards and Technology, Boulder, CO, United States

Correspondence*:

Jimmy Gammell

jimmy.gammell@colorado.edu

2 ABSTRACT

3 Equilibrium propagation is a learning framework that marks a step forward in the search for
4 a biologically-plausible implementation of deep learning, and is appealing for implementation
5 in neuromorphic analog hardware. However, previous implementations on layered networks
6 encountered a vanishing gradient problem that has not yet been solved in a simple, biologically-
7 plausible way. In this paper, we demonstrate that the vanishing gradient problem can be overcome
8 by replacing some of a layered network's connections with random layer-skipping connections.
9 This approach could be conveniently implemented in neuromorphic analog hardware, and is
10 biologically-plausible.

11 **Keywords:** equilibrium propagation, deep learning, small-world, layer-skipping connections, neuromorphic computing, biologically-
12 motivated

1 INTRODUCTION

13 Equilibrium propagation Scellier and Bengio [2016] is a learning framework for energy-based networks
14 such as the continuous Hopfield network [Hopfield, 1984]. It is appealing relative to backpropagation
15 because it is more biologically-plausible, and as a side-effect could be implemented more-easily in
16 neuromorphic analog hardware.

17 Implementation of equilibrium propagation in [Scellier and Bengio, 2016] was hindered by a vanishing
18 gradient problem whereby networks with as few as 3 hidden layers trained slowly on MNIST [LeCun and
19 Cortes, 1998] - a serious issue given that network depth is critical to performance on difficult datasets
20 [Simonyan and Zisserman, 2014; Srivastava et al., 2015b] and that convergence to a low error rate on
21 MNIST is a low bar to meet. The problem was overcome in [Scellier and Bengio, 2016] by independently
22 tuning a unique learning rate for each layer in the network; however, this approach is unappealing
23 because (1) it introduces additional hyperparameters to tune, (2) it would be inconvenient to implement in
24 neuromorphic analog hardware, and (3) it has not been observed in biological systems.

25 The purpose of this paper is to demonstrate that in this context the vanishing gradient problem can instead
26 be solved by randomly replacing some of a layered network's connections with layer-skipping connections.

¹ Through this modification we have achieved 0% training error (out of 50,000 examples) and $\lesssim 2.5\%$ test error (out of 10,000 examples) on MNIST using a network with three hidden layers and no regularization term in its cost function. These error rates are comparable to those of other biologically-motivated networks [Bartunov et al., 2018] and are roughly the same as those of the layered network with unique, manually-tuned learning rates in [Scellier and Bengio, 2016]. Our method could be implemented with relative ease in any system with configurable connectivity. Layer-skipping connections have been observed in biological brains [Bullmore and Sporns, 2009], so the approach is biologically-plausible. Similar techniques have seen success in convolutional [He et al., 2015; Srivastava et al., 2015a] and multilayer feedforward [Xiaohu et al., 2011; Krishnan et al., 2019] networks. Our findings outlined in this paper suggest that layer-skipping connections are effective-enough to be appealing in contexts where simplicity and biological plausibility are important.

2 METHODS

2.1 Equilibrium propagation

Similarly to backpropagation, equilibrium propagation [Scellier and Bengio, 2016] trains networks by approximating gradient descent on a cost function. Equilibrium propagation is applicable to any network with dynamics characterized by evolution to a fixed point of an associated energy function; our implementation is a recreation of that in [Scellier and Bengio, 2016], which applies it to a continuous Hopfield network [Hopfield, 1984]. The mathematical formulation of the framework can be found in [Scellier and Bengio, 2016].

A major reason backpropagation is not biologically-plausible is that to implement it, each neuron would need two distinct mechanisms for information transmission: one to transmit its activation to shallower neurons during the forward-propagation phase, and another to transmit error-correction information to deeper neurons during the backward-propagation phase [Bengio et al., 2015]. While this is easy in a digital computer that can oversee and manipulate an entire network, it would be cumbersome in hardware (biological or otherwise) consisting of many simple, independent computational nodes with limited ability to share information. In contrast, equilibrium propagation consists of a free phase (comparable to forward-propagation) and a weakly-clamped phase (comparable to backward-propagation) during which each neuron only needs to know the activations of neighboring neurons, so only one mechanism for information transmission is needed [Scellier and Bengio, 2016]. In a similar vein, to implement backpropagation each neuron would need mechanisms to compute both an activation value using activations of deeper neurons and error-correction information using that of shallower neurons. For equilibrium propagation a neuron would only need the ability to compute an activation given those of adjacent neurons [Scellier and Bengio, 2016].

2.1.1 Implementation in a continuous Hopfield network

Here we summarize the equations through which a continuous Hopfield network is trained using equilibrium propagation; this summary is based on the more-thorough and more-general treatment in [Scellier and Bengio, 2016].

Consider a network with n neurons organized into an input layer with p neurons, hidden layers with q neurons and an output layer with r neurons. Let the activations of these neurons be denoted respectively by vectors $\mathbf{x} \in \mathbb{R}^p$, $\mathbf{h} \in \mathbb{R}^q$ and $\mathbf{y} \in \mathbb{R}^r$, and let $\mathbf{s} = (\mathbf{h}^T, \mathbf{y}^T)^T \in \mathbb{R}^{q+r}$ and $\mathbf{u} = (\mathbf{x}^T, \mathbf{s}^T)^T \in \mathbb{R}^n$ be vectors of, respectively, the activations of non-fixed (non-input) neurons and of all neurons in the network.

¹ This modification was inspired by small-world topology [Watts and Strogatz, 1998]; however, we have not observed a strong correlation between network performance and common metrics of small-worldness (characteristic path length, clustering coefficient, small-world coefficient).

Let $\mathbf{W} \in \mathbb{R}^{n \times n}$ and $\mathbf{b} \in \mathbb{R}^n$ denote the network's weights and biases where w_{ij} is the connection weight between neurons i and j and b_i is the bias for neuron i ($\forall i$ $w_{ii} = 0$ to prevent self-connections), and let ρ denote its activation function; here and in [Scellier and Bengio, 2016],

$$\rho(x) = \begin{cases} 0 & x < 0 \\ x & 0 \leq x \leq 1 \\ 1 & x > 1 \end{cases} \quad (1)$$

is a hardened sigmoid function where $\rho'(0) = \rho'(1)$ is defined to be 1 to avoid neuron saturation. Let $\rho((x_1, \dots, x_n)^T) = (\rho(x_1), \dots, \rho(x_n))^T$.

The behavior of the network is to perform gradient descent on a total energy function F that is modified by a training example $(\mathbf{x}_d, \mathbf{y}_d)$. Consider energy function $E: \mathbb{R}^n \rightarrow \mathbb{R}$,

$$E(\mathbf{u}; \mathbf{W}, \mathbf{b}) = \frac{1}{2} \mathbf{u}^T \mathbf{u} - \frac{1}{2} \rho(\mathbf{u})^T \mathbf{W} \rho(\mathbf{u}) - \mathbf{b}^T \mathbf{u} \quad (2)$$

and arbitrary cost function $C: \mathbb{R}^r \rightarrow \mathbb{R}_+$; here and in [Scellier and Bengio, 2016] it is a quadratic cost function given by

$$C(\mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathbf{y}_d\|_2^2, \quad (3)$$

though the framework still works for cost functions incorporating a regularization term dependent on \mathbf{W} and \mathbf{b} . The total energy function $F: \mathbb{R}^n \rightarrow \mathbb{R}$ is given by

$$F(\mathbf{u}; \beta, \mathbf{W}, \mathbf{b}) = E(\mathbf{u}; \mathbf{W}, \mathbf{b}) + \beta C(\mathbf{y}) \quad (4)$$

where the clamping factor β is a small constant. \mathbf{s} evolves over time t as

$$\frac{d\mathbf{s}}{dt} \propto -\frac{\partial F}{\partial \mathbf{s}}. \quad (5)$$

Equilibrium has been reached when $\frac{\partial F}{\partial \mathbf{s}} \approx 0$. This can be viewed as solving the optimization problem

$$\underset{\mathbf{s} \in \mathbb{R}^{q+r}}{\text{minimize}} F((\mathbf{x}_d^T, \mathbf{s}^T)^T; \beta, \mathbf{W}, \mathbf{b}) \quad (6)$$

by using gradient descent to find a local minimum of F .

The procedure for training on a single input-output pair $(\mathbf{x}_d, \mathbf{y}_d)$ is as follows:

1. Clamp \mathbf{x} to \mathbf{x}_d and perform the free-phase evolution: evolve to equilibrium on the energy function $F(\mathbf{u}; 0, \mathbf{W}, \mathbf{b})$ in a manner dictated by equation 5. Record the equilibrium state \mathbf{u}^0 .
2. Perform the weakly-clamped evolution: evolve to equilibrium on the energy function $F(\mathbf{u}; \beta, \mathbf{W}, \mathbf{b})$ using \mathbf{u}^0 as a starting point. Record the equilibrium state \mathbf{u}^β .
3. Compute the correction to each weight in the network:

$$\Delta W_{ij} = \frac{1}{\beta} (\rho(u_i^\beta) \rho(u_j^\beta) - \rho(u_i^0) \rho(u_j^0)). \quad (7)$$

Adjust the weights using $W_{ij} \leftarrow W_{ij} + \alpha \Delta W_{ij}$ where the learning rate α is a positive constant.

88 4. Compute the correction to each bias in the network:

$$\Delta b_i = \frac{1}{\beta}(\rho(u_i^\beta) - \rho(u_i^0)) \quad (8)$$

89 and adjust the biases using $b_i \leftarrow b_i + \alpha \Delta b_i$.

90 This can be repeated on as many training examples as desired. Training can be done on batches by
 91 computing ΔW_{ij} and Δb_i for each input-output pair in the batch, and correcting using the averages of
 92 these values. Note that the correction to a weight is computed using only the activations of neurons it
 93 directly affects, and the correction to a bias is computed using only the activation of the neuron it directly
 94 affects. This contrasts with backpropagation, where to correct a weight or bias l layers from the output it is
 95 necessary to know the activations, derivatives and weights of all neurons between 0 and $l - 1$ layers from
 96 the output.

97 2.2 Vanishing gradient problem

98 Vanishing gradients are problematic because they reduce a network’s rate of training. and could be
 99 difficult to represent in neuromorphic analog hardware due to limited bit depth.

100 The vanishing gradient problem is familiar in the context of conventional feedforward networks, where
 101 techniques such as the weight initialization scheme in [Glorot and Bengio, 2010], the use of activation
 102 functions with derivatives that do not lead to output saturation [Schmidhuber, 2015], and batch norma-
 103 lization [Ioffe and Szegedy, 2015] have been effective at overcoming it. However, in the context of the
 104 networks trained in [Scellier and Bengio, 2016], the vanishing gradient problem persists even when the
 105 former two techniques are used. To our knowledge batch normalization has not been used in the context of
 106 equilibrium propagation; however, it seems unlikely to be biologically-plausible.

107 2.3 Implementation

We recreated the equilibrium propagation implementation in [Scellier and Bengio, 2016] using the Pytorch library.² Like the networks in [Scellier and Bengio, 2016], our networks are continuous Hopfield networks with a hardened sigmoid activation function

$$\sigma(x) = \text{Max}\{0, \text{Min}\{x, 1\}\}$$

and squared-error cost function with no regularization term

$$C = ||\mathbf{y} - \mathbf{y}_d||_2^2,$$

108 where \mathbf{y} is the network’s output and \mathbf{y}_d is the target output. Tests were run on MNIST [LeCun and Cortes,
 109 1998] grouped into batches of 20 examples, with the 50,000 training examples used for training and the
 110 10,000 validation examples used for computing test errors.

111 We use two performance-enhancing techniques that were used in [Scellier and Bengio, 2016]: we
 112 randomize the sign of β before training on each batch, which has a regularization effect, and we use
 113 persistent particles, where the state of the network after training on a given batch during epoch n is used as
 114 the initial state for that batch during epoch $n + 1$. Persistent particles reduce the computational resources
 115 needed to approximate the differential equation governing network evolution, and would be unnecessary
 116 in an analog implementation that can approximate the equation efficiently. Note that this technique leads

² https://github.com/jgammell/Equilibrium_Propagation_mobile.git

117 to higher error rates early in training than would be present with a more-thorough approximation of the
 118 differential equation.

119 2.3.1 Layered topology with per-layer rates

120 We recreated the 5-layer network evaluated in [Scellier and Bengio, 2016]. It has the standard layered
 121 topology shown in 1, and consists of a 784-neuron input layer, 3 500-neuron hidden layers and a 10-neuron
 122 output layer. Weights are initialized using the scheme from [Glorot and Bengio, 2010]. As mentioned
 123 above, each layer has a unique learning rate; the rates are $\alpha_1 = .128$, $\alpha_2 = .032$, $\alpha_3 = .008$ and $\alpha_4 = .002$
 124 where α_i is the learning rate for the connection weights between layers i and $i + 1$ and for the biases in
 125 layer i , and the input and output layers are denoted $i = 1$ and $i = 5$, respectively.

126 2.3.2 Layered topology with global learning rate

127 To illustrate the vanishing gradient problem and provide a point of reference, we also tested the network
 128 in section 2.3.1 with a single global learning rate of .02.

129 2.3.3 Our topology

Algorithm 1: Algorithm to produce our topology

Input: Layered network from section 2.3.2

Input: Integer n , giving number of connections to replace

Output: A network with our modified topology

for hidden layer in network **do**

 | Add edge between each pair of neurons in layer

for $i \leftarrow 1$ **to** n **do**

 Randomly select pre-existing connection in network;

 Add connection between random unconnected pair of
 neurons in network;

 // Do not allow self connections

 // Do not allow connections between two input neurons or
 between two output neurons

 Remove pre-existing connection;

return modified network

130 To generate a network with our topology, we use algorithm 1. This topology is illustrated in figure 2.
 131 The above algorithm is approximately equivalent to the algorithm for generating a small-world network
 132 described in [Watts and Strogatz, 1998] with $p = 1 - (\frac{N_o - 1}{N_o})^n$ for $p \lesssim .2$, where N_o is the number
 133 of connections in the network; to contextualize the number of replaced connections we will henceforth
 134 describe networks with our topology in terms of p instead of n . We have seen good results with $p \approx 8\%$.
 135 We have seen similar results when connections are added to the network, rather than randomly replaced
 136 (algorithm 1, without removing pre-existing connections).

137 For these networks we use a global learning rate of .02 and, as in the networks from sections 2.3.1 and
 138 2.3.2, initialize connections between neurons in adjacent layers using the scheme from [Glorot and Bengio,
 139 2010]. For all other connections we draw initial weights from the uniform distribution $U[-.05, .05]$ where
 140 the value .05 was determined empirically to yield good results.

3 RESULTS

141 We compared the networks described in section 2.3 by observing their behavior while training on MNIST
 142 [LeCun and Cortes, 1998]. All networks used $\epsilon = .5$, $\beta = 1.0$, 500 free-phase iterations, 8 weakly-clamped-
 143 phase iterations, and were trained for 250 epochs.

3.1 Network performance comparison

Figure 3 illustrates that our network significantly outperforms one with a global learning rate, and achieves close to the same training and test error rates as one with unique learning rates, albeit after around 25% more epochs. Both our network and the layered network with unique learning rates achieve approximately a 2.5% test error and 0% training error, whereas the layered network with a global learning

3.2 Training rates of individual pairs of layers

rate has test and training error rates around .5% higher than the other two networks; it is unclear whether it would converge given enough time, but it is clear that it is inferior. To observe the extent of the vanishing gradient problem, for each network we tracked the root-mean-square correction to weights in each of its layers during training on MNIST [LeCun and Cortes, 1998].

Figure 4 shows an 11-point centered moving average of these values (without averaging the values are very volatile). It can be seen that for the layered network with a global learning rate, the magnitude of the correction to a typical neuron vanishes with depth relative to the output, with the shallowest weights training around 100 times faster than the deepest weights - this illustrates the vanishing gradient problem. The use of unique learning rates is very effective at making corrections uniform. Our topology with $p = 7.56\%$ is effective at making deeper layers train in a uniform way, but the output layer still trains around 10 times faster than deeper layers; nonetheless, figure 3 suggests that this imperfect solution still yields a significant performance benefit.

The fast training of the output layer in the network with our topology is probably because no layer-skipping connections attach directly to the target output, so for any value of p the shortest path between a deep neuron and the target layer is at least 2 connections long, whereas the path between an output neuron and the target layer is only 1 connection long.

We tracked the training error after one epoch of a network with our topology while varying p ; the results are shown in figure 5. For $p < .1\%$, there is little improvement in the error rate as p is increased, but there is substantial improvement in the uniformity of the training rates of deep layers. When $p > .1\%$, the deep layers are very uniform, and the error rate starts decreasing with p at a rate that is slightly slower than exponential; at this point there is little improvement in the uniformity of deep layers, but the rate of the shallowest layer appears to move closer to those of the deeper layers.

We found that our topology performs significantly worse than the basic topology with one learning rate when few connections are replaced. This could be due to a poor weight initialization scheme for the added intralayer connections; we have noticed anecdotally that networks appear to be less-sensitive to their weight initialization scheme as connections are replaced. We have found that networks perform poorly relative to a basic network with one learning rate until p is in the ballpark of 7%. This experiment suggests that training rate will keep improving long after that, but does not show long-term performance or test performance; we suspect that a network's generalization ability will suffer for large p as it loses its regimented nature.

4 DISCUSSION

4.1 Comparing the computational complexity of equilibrium propagation and backpropagation

4.1.1 Requirements of equilibrium propagation

It follows from equations 2, 4 and 5 that to determine its state, the i -th neuron in a network must compute

$$\frac{\partial F}{\partial u_i} = u_i - \frac{1}{2} \rho'(u_i) \left[\sum_{j \neq i} W_{ij} \rho(u_j) + b_i \right],$$

plus the term $\beta(u_i - y_i^{target})$ for output neurons when using a squared-error cost function, and then integrate the result over time. Parameter correction rules are given by equations 7 and 8. A block diagram of this process is shown in figure 6 and qualitatively describes one way equilibrium propagation could potentially be implemented in hardware.

4.1.2 Requirements of backpropagation

In backpropagation, the activation value of a neuron i in layer l is given by

$$\rho(u_i^l) = \rho\left(\sum_j W_{ij}^l u_j^{l-1} + b_i^l\right).$$

Parameters are then updated by computing error correction terms δ_i^l for each neuron i in layer l ; for the output layer L the correction is

$$\delta_i^L = \rho'(u_i^L) (\rho(u_i^L) - y_i^{target})$$

and for deeper layers it is

$$\delta_i^l = \rho'(u_i^l) \sum_j W_{ij}^{l+1} \delta_j^{l+1}.$$

Weights are corrected using

$$\Delta W_{ij}^l = \rho(u_i^{l-1}) \delta_j^l$$

and biases using

$$\Delta b_i^l = \delta_i^l.$$

The block diagram in figure 7 qualitatively describes a way this algorithm could potentially be implemented in hardware.

4.1.3 Comparison

The most-significant difference between the algorithms is that in equilibrium propagation, the free and weakly-clamped phases of training are identical for most neurons and the weakly-clamped phase requires only slight modification to output neurons, whereas in backpropagation these phases very different functionality from essentially all neurons. Another visible difference is that in equilibrium propagation each neuron corresponds to a single synapse whereas in backpropagation a neuron corresponds to two synapses; we do not expect this difference to be significant because a synapse in the former case a synapse takes as inputs the outputs of all neighboring neurons, whereas in the latter case each has inputs from either shallower or deeper neurons (about half as many). While neurons in equilibrium propagation explicitly write their states to memory after the free phase, in backpropagation the need for distinct state variables to hold the activation and error term of each neurons implies a need for the same amount of memory. Various characteristics of both algorithms are compared side-by-side in table 1.

4.2 Related work

4.3 Directions for Future Research

There are several directions in which future research could be taken:

- Evaluating the effectiveness of this approach on hard datasets, such as CIFAR and ImageNet.
- Evaluating the effect of p on a network's test error.
- Exploring the effectiveness of layer-skipping connections on deeper networks.
- Exploring the effectiveness of a network when layer-skipping connections are used during training and removed afterwards.

CONFLICT OF INTEREST STATEMENT

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

REFERENCES

- Bartunov, S., Santoro, A., Richards, B. A., Hinton, G. E., and Lillicrap, T. P. (2018). Assessing the scalability of biologically-motivated deep learning algorithms and architectures. *CoRR* abs/1807.04587
- Bengio, Y., Lee, D., Bornschein, J., and Lin, Z. (2015). Towards biologically plausible deep learning. *CoRR* abs/1502.04156
- Bullmore, E. and Sporns, O. (2009). Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature*
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, eds. Y. W. Teh and M. Titterton (Chia Laguna Resort, Sardinia, Italy: PMLR), vol. 9 of *Proceedings of Machine Learning Research*, 249–256
- He, K., Zhang, X., Ren, S., and Sun, J. (2015). Deep residual learning for image recognition. *CoRR* abs/1512.03385
- Hopfield, J. (1984). Neurons with graded response have collective computational properties like those of two-state neurons. *Proceedings of the National Academy of Sciences of the United States of America* 81, 3088–92. doi:10.1073/pnas.81.10.3088
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR* abs/1502.03167
- [Dataset] Krishnan, G., Du, X., and Cao, Y. (2019). Structural pruning in deep neural networks: A small-world approach
- [Dataset] LeCun, Y. and Cortes, C. (1998). The mnist database of handwritten digits
- [Dataset] Scellier, B. and Bengio, Y. (2016). Equilibrium propagation: Bridging the gap between energy-based models and backpropagation
- Schmidhuber, J. (2015). Deep learning in neural networks: An overview. *Neural Networks* 61, 85–117. doi:10.1016/j.neunet.2014.09.003
- [Dataset] Simonyan, K. and Zisserman, A. (2014). Very deep convolutional networks for large-scale image recognition
- Srivastava, R. K., Greff, K., and Schmidhuber, J. (2015a). Highway networks. *CoRR* abs/1505.00387
- [Dataset] Srivastava, R. K., Greff, K., and Schmidhuber, J. (2015b). Training very deep networks
- Watts, D. and Strogatz, S. (1998). Collective dynamics of 'small-world' networks. *Nature*

241 Xiaohu, L., Xiaoling, L., Jinhua, Z., Yulin, Z., and Maolin, L. (2011). A new multilayer feedforward small-
242 world neural network with its performances on function approximation. In *2011 IEEE International*
243 *Conference on Computer Science and Automation Engineering*

FIGURES

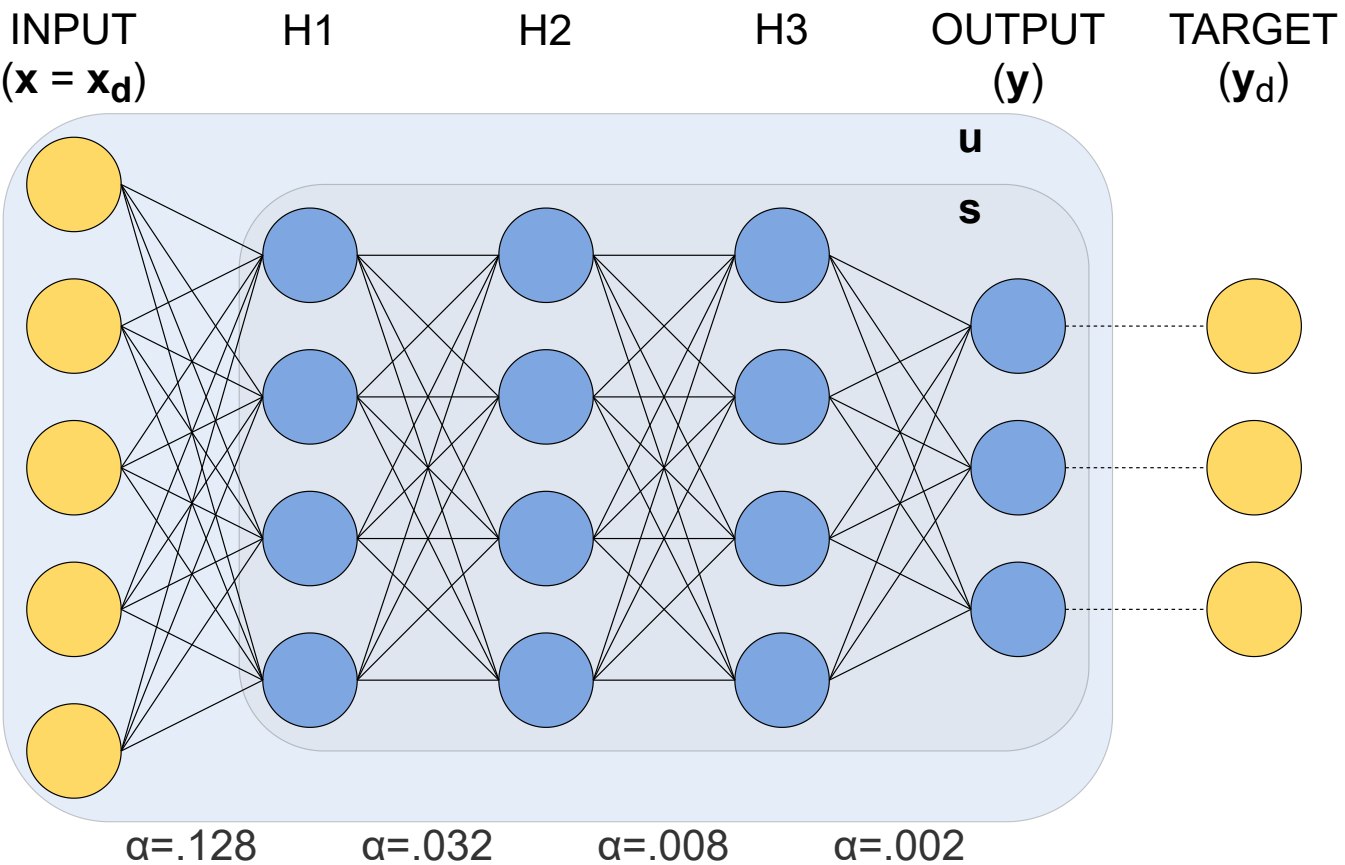


Figure 1. Topology of the layered network tested in [Scellier and Bengio, 2016]. All pairs of neurons in adjacent layers are connected. All connections are bidirectional. To compensate for the vanishing gradient problem, the learning rate is reduced by a factor of 4 each time distance from the output decreases by one layer.

TABLES

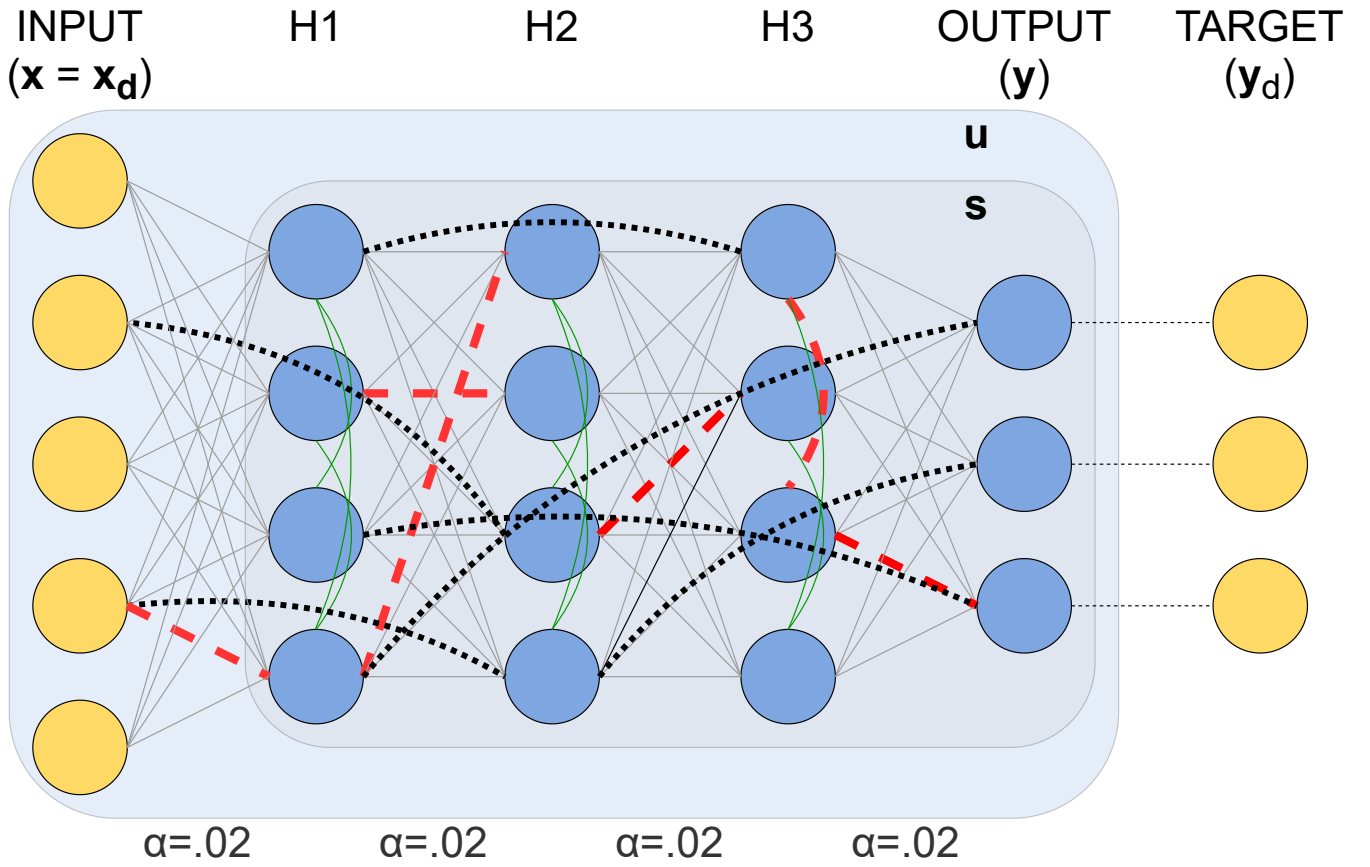


Figure 2. Our modifications to the topology of figure 1 to avoid a vanishing gradient while using a global learning rate. Red dotted lines denote connections that have been removed, black dotted lines denote their replacements, and green solid lines denote added intralayer connections. All connections are bidirectional. This illustration shows a network with $p = 8\%$.

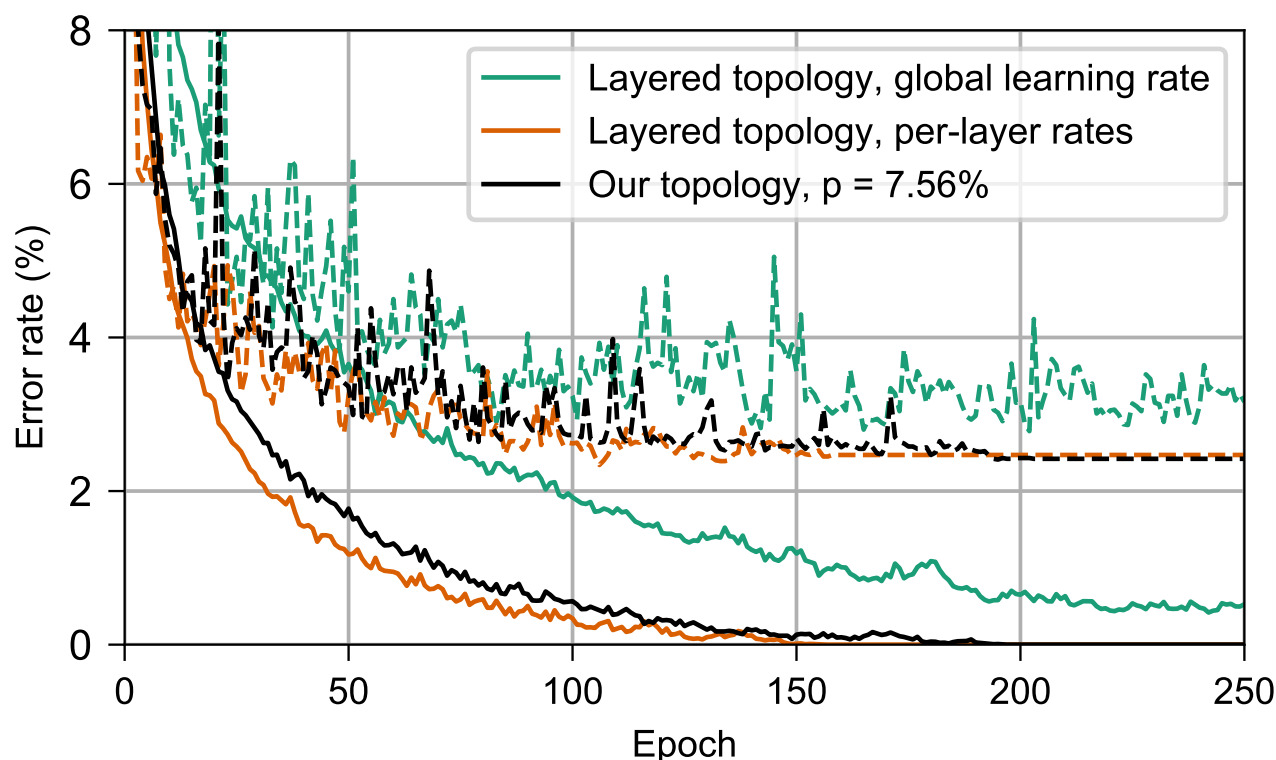


Figure 3. Performance on MNIST of the networks in section 2.3. Dashed lines show the test error and solid lines show the training error. In green is a layered network with a global learning rate (section 2.3.2), in orange is a layered network with per-layer rates individually tuned to counter the vanishing gradient problem (section 2.3.1), and in green is a network with our topology, $p = 7.56\%$ (section 2.3.3). Observe that our topology is almost as effective as per-layer rates at countering the vanishing gradient problem that impedes training of the layered network with a global learning rate.

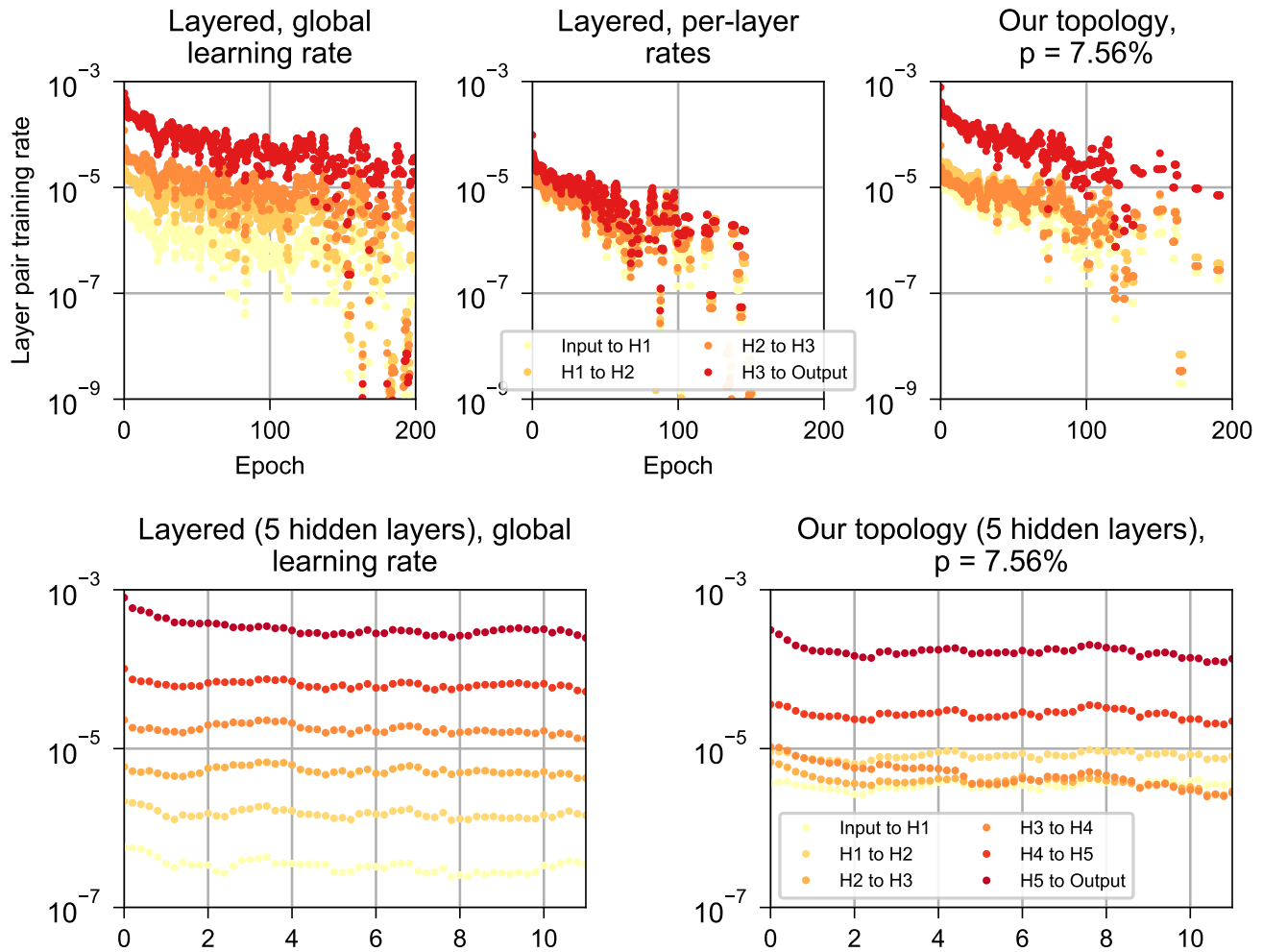


Figure 4. Root-mean-square corrections to weights in different layers while training on MNIST, for the networks in section 2.3. For clarity, values were subjected to an 11-point centered moving average. (top left) A layered network with a single global learning rate (section 2.3.2). (top center) A layered network with a unique, individually-tuned learning rate for each layer (section 2.3.1). (top right) A network with our topology, $p = 7.56\%$ (section 2.3.3). (bottom left) A layered network with 5 100-neuron hidden layers and a single global learning rate (section ??). (bottom right) A network with our topology, $p = 7.56\%$, and 5 100-neuron hidden layers (section ??). Observe that for the shallower networks the layered topology with a global learning rate has a vanishing gradient problem, which is almost completely solved by tuning an individual learning rate for each layer. Our topology improves the situation by making training uniform among the deeper layers, although the shallowest layer still trains more-quickly than the deeper layers. For the deeper networks, the same trend is apparent but not as strong; we believe this is due to sub-optimal hyperparameter settings.

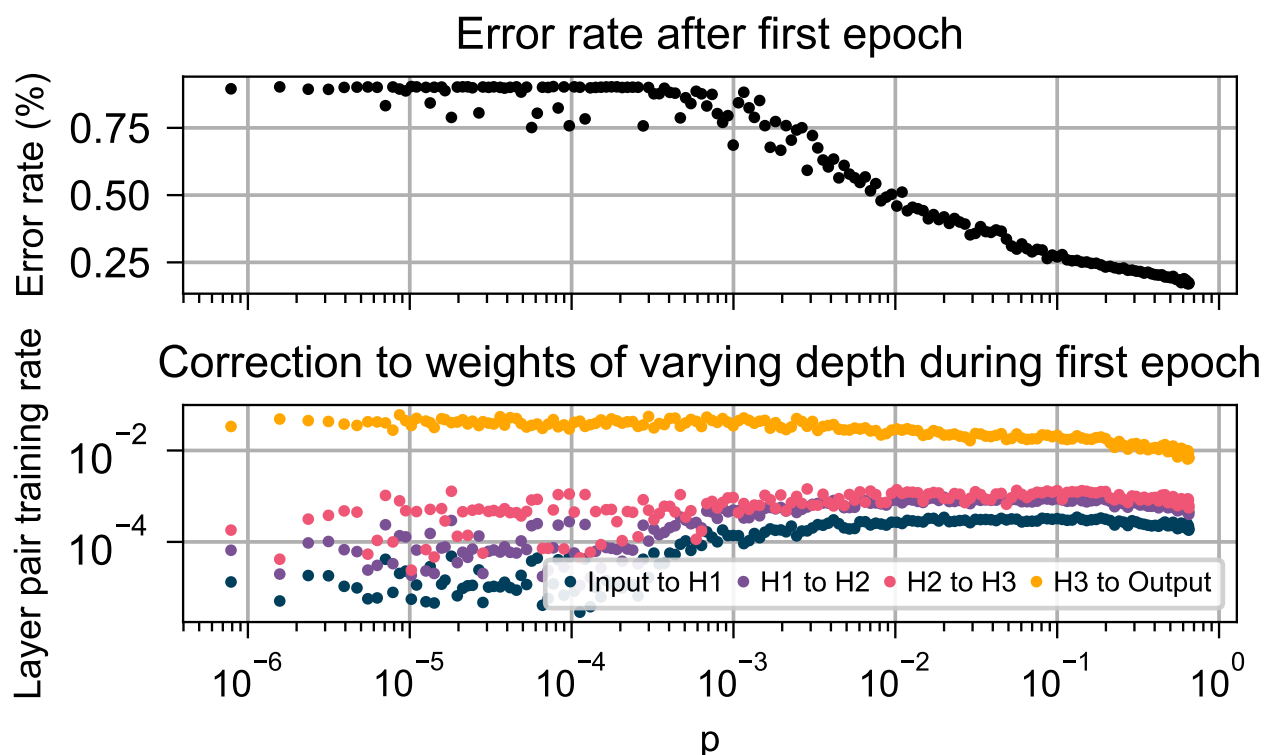


Figure 5. Behavior of our network (section 2.3.3) with varying p , during the first epoch of training. (top) The training error after one epoch. (bottom) Root-mean-square correction to weights in different layers during the first epoch. Observe that as p is increased, the error rate decreases and the root-mean-square corrections to each layer become more-uniform.

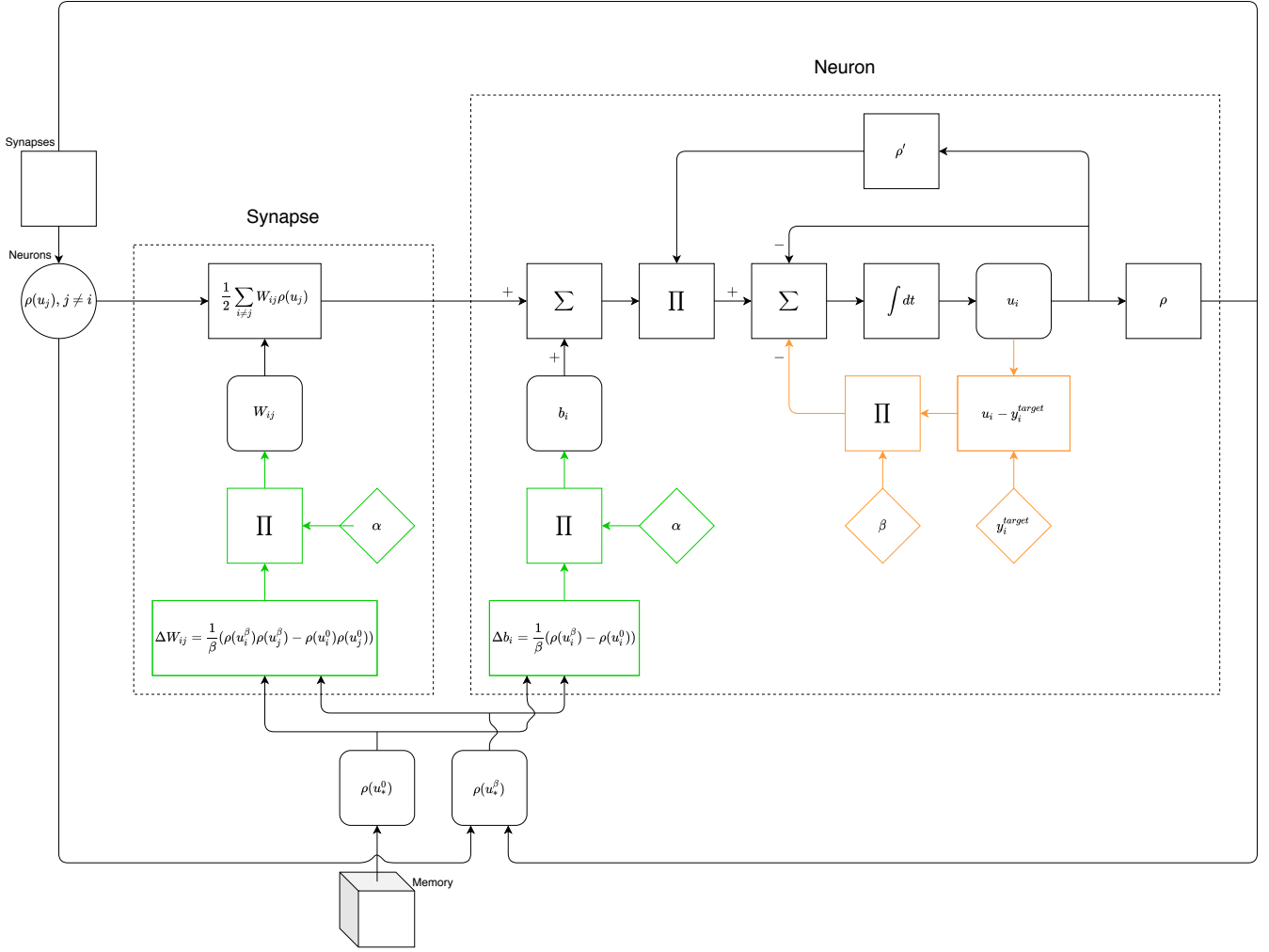


Figure 6. Illustration of the functionality needed to implement equilibrium propagation in hardware. Black lines denote functionality needed in the free phase. Green lines denote functionality to correct parameters. Orange lines denote functionality needed only by output neurons, that is unique to the weakly-clamped phase. There is no functionality unique to the weakly-clamped phase that is needed by all neurons.

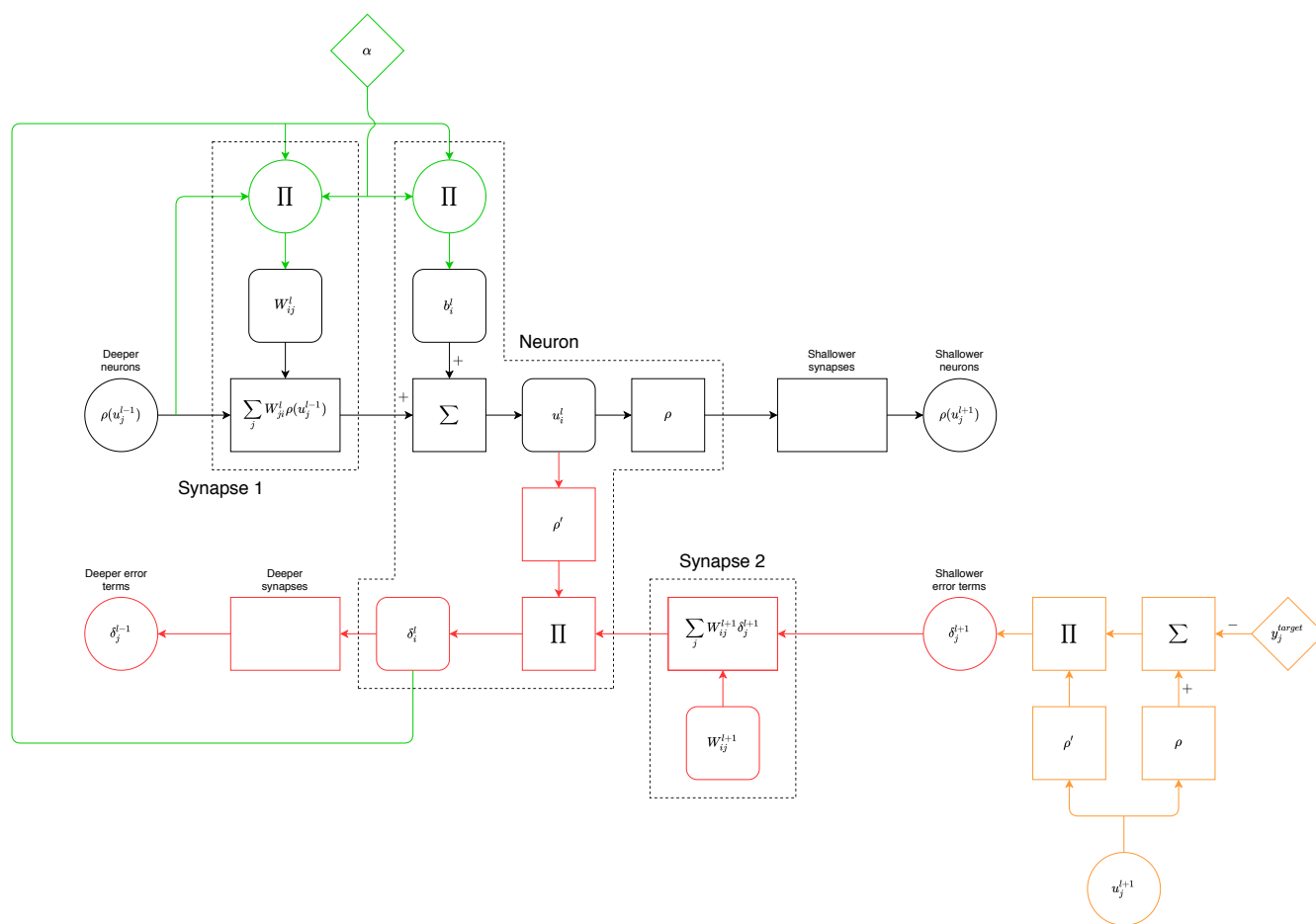


Figure 7. Illustration of the functionality needed to implement backpropagation in hardware. Black lines denote functionality needed in the forwards phase. Green lines denote functionality to correct parameters. Red lines denote functionality unique to the backwards phase that is needed by all neurons. Orange lines denote functionality needed only by output neurons, that is unique to the backwards phase.

	Backpropagation	Equilibrium Propagation
Memory	Space to store activation and error term for each neuron	Space to store free and weakly-clamped activations for each neuron
Nonlinear activation function	Yes	Yes
Derivative of nonlinear activation function	Yes	Yes
Number of distinct computations	2 - computations during forwards and backwards phases are distinct	≈ 1 - hidden neurons perform same computation in both phases. Output neurons perform a similar but modified version of the same computation.
Types of connections	Unidirectional to transmit activation to shallower neighbors and error to deeper neighbors	Bidirectional to each neighbor
Correction computation	Corrections require dedicated circuitry unique from that implementing propagation	Corrections require dedicated circuitry unique from that implementing evolution
Order of computations	Forwards propagation phase where layers are computed from deepest to shallowest; backwards propagation phase where layers are computed from shallowest to deepest; parameter update phase	Free phase where all neurons evolve simultaneously; weakly-clamped phase where all neurons evolve simultaneously; parameter update phase

Table 1.