

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

Jimmy Gammell
Sae Woo Nam
Adam McCaughan

ABSTRACT

Equilibrium propagation is a learning framework for energy-based networks that can be implemented by neurons that perform only one type of computation in the prediction and correction phases of training, and that computes parameter corrections for a given neuron using only the activation values of directly-connected neurons. This makes it an appealing candidate for implementation in neuromorphic analog hardware, and marks a step forward in the search for a biologically-plausible implementation of deep learning. However, in previous implementations of equilibrium propagation, layered networks suffered from a vanishing gradient problem that has not yet been solved in a simple, biologically-plausible way. In this paper, we demonstrate that the vanishing gradient problem can be counteracted by replacing some of a layered network's connections with random layer-skipping connections in a manner inspired by small-world networks. This approach could be conveniently implemented in neuromorphic analog hardware and is biologically-plausible.

CCS CONCEPTS

• Computing methodologies → Bio-inspired approaches; Neural networks.

KEYWORDS

deep learning, neural network, biologically-motivated, equilibrium propagation, vanishing gradient, small-world, neuromorphic hardware

ACM Reference Format:

Jimmy Gammell, Sae Woo Nam, and Adam McCaughan. 2020. Layer-skipping connections facilitate training of layered networks using equilibrium propagation. In *Proceedings of International Conference on Neuromorphic Systems (ICONS 2020)*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

The equilibrium propagation learning framework [17] is a method for training a class of energy-based networks, the prototype for which is the continuous Hopfield network [7]. It is appealing as

a framework that could be implemented in neuromorphic analog hardware because unlike in backpropagation, neurons are required to perform only one type of computation in the prediction and correction phases of training, and the parameters of a neuron can be updated using only the activations of neurons to which it is directly connected - there is no need for feedback paths through which to transmit information about the parameters and states of neurons throughout the entire network. Backpropagation is not biologically plausible, and a major reason is that credit assignment to a given neuron requires precise knowledge of the nonlinearities and derivatives of all neurons in the feedforward path between that neuron and the output[2]; equilibrium propagation avoids this issue. Another major reason is that it is implausible for the forward and backwards propagation phases to use different computations requiring different computational circuits [17]; equilibrium propagation also avoids this issue.

It has been demonstrated [17] that a continuous Hopfield network with a basic multilayer topology can be trained on MNIST [11] through the equilibrium propagation framework. However, previous implementations encountered a vanishing gradient problem that significantly impedes training of networks with several hidden layers. Given that network depth is critical for performance on difficult datasets [20, 21], and that convergence to a low error rate on MNIST is a low bar for a network to meet, this is a nontrivial issue. It has been demonstrated [17] that the problem can be solved by using a unique learning rate for parameters at different depths in the network, with deeper parameters trained with larger learning rates to counteract gradient attenuation with depth. This approach is unappealing because (1) it introduces additional hyperparameters that must be tuned, (2) it would be inconvenient to implement and tune unique learning rates in analog hardware, and (3) it has not been observed in biological systems.

The purpose of this paper is to introduce a modification to the basic multilayer topology that can counteract the vanishing gradient problem in the context of energy-based networks trained using equilibrium propagation. The modification is inspired by small-world topology [23], and consists of first fully-connecting a network's layers, then replacing some of its connections with random layer-skipping connections. We achieve 0% training error and under 2.5% test error on MNIST using a network with 3 hidden layers and no regularization term in its cost function. These error rates are comparable to those achieved by other biologically-motivated networks [1] and are the same as those achieved by a basic multilayer network with 3 hidden layers in the original paper using a basic topology and manually-tuned per-layer learning rates [17], albeit they are achieved after around 25% more epochs. Our method adds only one additional hyperparameter - the number of connections to randomly replace - and could be implemented

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICONS 2020, July 2020, Chicago, Illinois

© 2020 Association for Computing Machinery.

ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

with relative ease in any system with configurable connectivity.¹ Layer-skipping connections have been observed in biological brains [3], so this approach is biologically-plausible. Similar techniques have been used with some success in convolutional [6, 22] and basic multilayer feedforward [10, 24] networks with varying degrees of success. Our findings in this paper show that layer-skipping connections are effective enough to be appealing in contexts where simplicity and biological plausibility are important.

2 BACKGROUND AND THEORY

2.1 Equilibrium propagation

Equilibrium propagation [17] is a learning framework for energy-based networks that trains their parameters by approximating gradient descent on some arbitrary cost function. It is applicable to any network with dynamics characterized by evolution to a fixed point of an associated energy function; in this and in [17], it is implemented on a continuous Hopfield network [7]. For a variety of reasons outlined in [2], backpropagation is not biologically-plausible. One of the major reasons is that to correct the parameters of a given neuron, backpropagation requires precise information about the activations and nonlinearities of all neurons in the corresponding feedforward path; equilibrium propagation avoids this problem by correcting a neuron's parameters using only the activations of neurons to which it directly connects. Another major reason backpropagation is not biologically-plausible is that it requires neurons to perform distinct computations (implemented using distinct circuitry) in the forward and backward propagation phases of training; in contrast, equilibrium propagation requires only one behavior of neurons in both the prediction (free) and correction (weakly-clamped) phases of training. In addition to enhancing its biological plausibility, these traits make equilibrium propagation appealing as a framework to implement on neuromorphic analog hardware because they would limit the complexity required of neurons and the amount of infrastructure needed to operate and train them.

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 [17].

Consider a network consisting of input, output and hidden layers with activations denoted respectively by \mathbf{x} , \mathbf{y} and \mathbf{h} , and let $\mathbf{s} = \{\mathbf{h}, \mathbf{y}\}$ and $\mathbf{u} = \{\mathbf{x}, \mathbf{h}, \mathbf{y}\}$. Let \mathbf{W} and \mathbf{b} denote the network's weights and biases and ρ its activation function; here and in [17]

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

is a hardened sigmoid function with

$$\rho'(x) := \begin{cases} 0 & x < 0 \\ 1 & 0 \leq x \leq 1 \\ 0 & x > 1 \end{cases} \quad (2)$$

¹Our current implementation also requires manual tuning of the initial weights of intralayer and layer-skipping connections, but we have found that networks are not very sensitive to these initializations and that performance is good as long as they are in the same ballpark as the rest of the network's initializations.

where ρ' is defined to be 1 at endpoints of its support to avoid saturation.

The behavior of the network is to perform gradient descent on a total energy function F . Consider energy function

$$E(\mathbf{u}) = \frac{1}{2} \sum_i u_i^2 - \frac{1}{2} \sum_{i \neq j} W_{ij} \rho(u_i) \rho(u_j) - \sum_i b_i \rho(u_i). \quad (3)$$

and arbitrary cost function $C(\mathbf{u}, \mathbf{W}, \mathbf{b}, \mathbf{y}_d)$; here and in [17] it is a quadratic cost function given by

$$C(\mathbf{y}, \mathbf{y}_d) = \frac{1}{2} \|\mathbf{y} - \mathbf{y}_d\|_2^2. \quad (4)$$

The total energy function is given by

$$F(\mathbf{u}, \mathbf{y}_d, \beta) = E(\mathbf{u}) + \beta C(\mathbf{y}, \mathbf{y}_d) \quad (5)$$

where the clamping factor β is a small constant. Gradient descent on F is defined by an equation of motion

$$\frac{ds}{dt} = -\frac{\partial F}{\partial s}; \quad (6)$$

equilibrium has been reached when the network is at a local minimum of F .

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

- (1) Perform the free-phase evolution: evolve to equilibrium on the energy function $F(\mathbf{u}, \mathbf{y}_d, 0)$ in a manner dictated by equation 6. Record the equilibrium state \mathbf{u}^0 .
- (2) Perform the weakly-clamped evolution: evolve to equilibrium on the energy function $F(\mathbf{u}, \mathbf{y}_d, \beta)$ 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.

- (4) Compute the correction to each bias in the network:

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

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

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

2.1.2 Approximation of equation of motion. In analog hardware the dynamics described by equation 6 could be implemented efficiently using leaky integrator neurons, but on digital hardware it is necessary to discretize and approximate these dynamics; we now describe the approximation used here and in [17]. Let $\mathbf{s}[n]$ denote the state of the network after n iterations of the approximation, N denote the total number of iterations and ϵ the size of each iteration. $\mathbf{s}[0]$ is the initial state of the network; here and in [17], $\mathbf{s}[0] = \mathbf{0}$ before the first epoch, and in subsequent epochs $\mathbf{s}[0]$ is

the equilibrium state of the network from the last time it executed a free phase on the same batch. Let

$$s_i[n] = s_i[n-1] - \epsilon \frac{\partial F}{\partial s_i}(\mathbf{u}, \mathbf{y}_d, \beta), \quad n = 1, \dots, N \quad (9)$$

(where $\beta = 0$ in the free phase); then $s[N]$ approximates the state after time ϵN that would be given by equation 6. We denote the number of iterations in the free and weakly-clamped phases N_{free} and $N_{weakly-clamped}$, respectively. It is only necessary to run the weakly-clamped phase long enough to observe the initial direction in which the network evolves and to allow the perturbation of the output to influence all layers in the network, so typically for a network with L layers, $N_{free} \gg N_{weakly-clamped} > L$.

2.2 Vanishing gradient problem

It has been observed [17] that energy-based networks with a layered topology trained using equilibrium propagation suffer from a vanishing gradient problem whereby the rate of training of a weight decreases approximately exponentially with the number of layers between that weight and the output layer. This is problematic because it dramatically increases the time needed to train a network, and gradients spanning several orders of magnitude could be hard to represent in analog implementations of the framework with limited bit depth. The vanishing gradient problem is familiar in the context of conventional networks trained through backpropagation, where the composition of many layers can impede training by amplifying or attenuating the effect of changes to deep neurons as these changes propagate to the output layer, leading to a vanishing or exploding gradient. For those networks the problem can be effectively addressed by initializing weights to make activation variances and backpropagated gradient variances approximately uniform with respect to depth as described in [5], and by using activation functions with derivatives that do not cause neurons to saturate [18], such as rectified linear units $\rho(x) = \text{Max}\{0, x\}$ or hardened sigmoids (equation 1). Batch normalization [9] is another effective technique for promoting uniform training in a network; inputs to each layer are whitened (the input distribution \mathbf{x} to each layer is transformed to have $E[\mathbf{x}] = 0$ and $\text{Var}[\mathbf{x}] = 1$) during each batch to prevent changes to deep parameters from changing the input distribution to shallower layers, requiring adjustment to compensate.

The weight initialization scheme from [5] and the use of hardened sigmoid activation functions do not solve the vanishing gradient problem in the context of equilibrium propagation. Batch normalization has not been tried in this context, but is not biologically-plausible. That the former two methods are ineffective suggests that the cause of the vanishing gradient problem is different than that in a conventional network. The situation in energy-based networks is more-complicated because there is no straightforward causal relationship between the activations of two given neurons; instead, neurons evolve as a coupled dynamical system.

In [17] the vanishing gradient problem was solved by using unique learning rates for each layer, chosen to make the magnitudes of corrections to neurons' parameters uniform regardless of depth in the network. While this method is effective, it is unappealing because it seems unlikely to take place in a biological brain, it would add complexity to an analog implementation of the

framework, and it introduces more hyperparameters that must be tuned to train a network. Our topology solves the problem without these unappealing characteristics.

A possible contributor to the vanishing gradient problem is the nature of the discrete approximation of a network's equation of motion; in appendix B we describe this issue and why we do not believe it to be the primary cause.

2.3 Small-world networks

Our topology was inspired by small-world graph topology as described in [23]. We have observed that energy-based networks with layered topology experience exponential attenuation of the gradient of a given parameter as its depth in the network increases, so it seems reasonable to expect that the vanishing gradient problem can be reduced by decreasing the number of connections needed to move between a given neuron and the output layer. Small-world topology offers a way to do this while largely preserving clustering between neurons within the same layer. It is also a topology that has been observed in biological brains [3]. Essentially, a layered network with fully-connected layers can be viewed as a network with a large typical path length between pairs of neurons and a large amount of clustering between nearby neurons. By randomly replacing a small proportion of its connections, its typical path length can be greatly reduced without significantly reducing its clustering. Quantitative metrics of a network's small-worldness are described in appendix A, but in our experiments these metrics have little correlation with network behavior (we discuss this in more detail later).

2.3.1 Algorithm for generating a small-world network. The following algorithm was introduced in [23] to convert a regimented graph into a small-world network. It was implemented on a lattice network but works equally-well for a layered neural network.

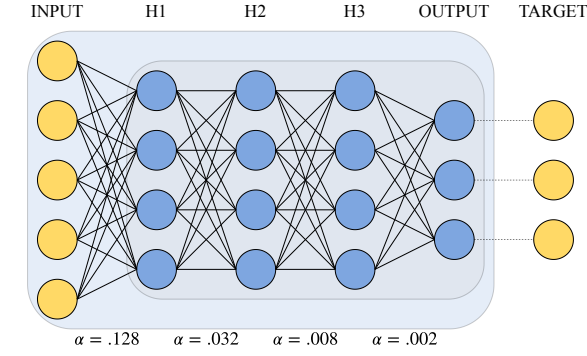
- (1) Generate a regimented graph.
- (2) Consider each preexisting edge in the graph, and with probability p :
 - (a) Randomly select a pair of vertices that are not connected and add an edge between them.
 - (b) Remove the preexisting edge.

As $p \approx 0$ is increased, the mean shortest path between a pair of vertices in the network (L from appendix A) decreases rapidly while preexisting clusters remain largely intact (C from appendix A does not significantly decrease). To generate a small-world network, p is chosen so that the network has a large amount of clustering and a small mean shortest path, relative to a random graph with the same number of nodes and edges (σ from appendix A is large).

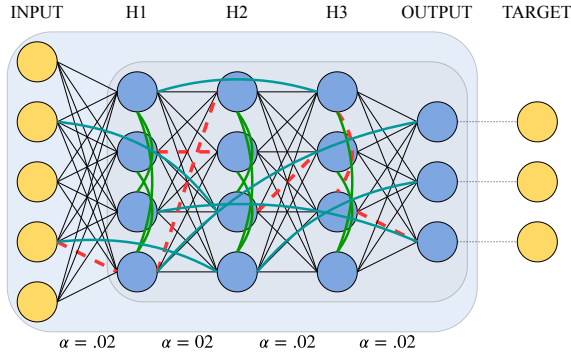
3 IMPLEMENTATION

We implemented ² the equilibrium propagation learning framework [17] using the Pytorch library, and verified that we could recreate the experiments run in [17]. We follow the procedure described in section 2.1; we use a hardened sigmoid activation function (equation 1) and a squared-error cost function with no regularization term (equation 4). Testing was done on the MNIST dataset [11] with input-output pairs grouped into batches of 20; training was done

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



(a) Topology of the basic multilayer network tested in [17]. All pairs of neurons in adjacent layers are connected, and there are no additional connections. The learning rate for weights is reduced by a factor of 4 each time distance from the output decreases by one layer, to compensate for the vanishing gradient problem.



(b) Changes we have made to the basic topology to compensate for the vanishing gradient problem while using a single learning rate for all weights. Red dotted lines denote connections that have been removed and blue lines denote their replacements. Green lines denote added connections within layers (these are also candidates for replacement). In this illustration layers have been made fully connected, and each connection has then been replaced by a random layer-skipping connection with probability $p \approx 8\%$.

Figure 1

on the 50,000 training pairs and testing on the 10,000 validation pairs.

We use two performance-boosting techniques described in [17]. We randomize the sign of β at the outset of each batch, which provides a regularization effect. We use persistent particles: after the first epoch, before training on a batch the network's state is initialized to its free equilibrium state from the last time it trained on the same batch. This reduces computational resources needed for training by reducing the number of iterations of equation 6 needed to reach equilibrium, and would be unnecessary in an analog implementation; note that training and test error rates during early epochs are artificially inflated as a result of this technique.

3.1 Basic topology with unique learning rates

We recreated the 5-layer network used in [17]. Its topology is illustrated in figure 1a.

The network consists of an input layer with 784 neurons, 3 hidden layers with 500 neurons each, and an output layer with 10 neurons. There are connections between every pair of neurons in adjacent layers, no connections between neurons in the same layer, and no layer-skipping connections.

The weight matrix was initialized using the Glorot-Bengio initialization scheme [5]: a weight connecting layer j with n_j neurons to layer $j + 1$ with n_{j+1} neurons is drawn from a uniform distribution

$$U\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}\right]. \quad (10)$$

For effective training it is necessary to use unique learning rates to train the weights connecting each pair of layers. We use the same learning rates that were used for experiments in [17]: .128, .032, .008, .002, in order of the deepest to the shallowest layer; to counter the vanishing gradient problem the learning rate is reduced by a factor of 4 each time the distance from the output layer decreases by one layer. Each neuron has a bias term that is trained using the learning rate corresponding to the weights of the neuron's input connections, e.g. the biases of the first hidden layer are trained with a learning rate of .128.

3.2 Basic topology with single learning rate

To provide a point of reference, we ran tests on a network that is identical to that in section 3.1 except with a single learning rate of .02 across the entire network. It performs poorly due to the vanishing gradient problem.

3.3 Our topology

To generate a network with our topology, we use the following procedure:

- (1) Generate a network with the topology described in section 3.1.
- (2) Add intralayer connections: Within each hidden layer (not within the input or output layers), add a connection between each pair of neurons.
- (3) Replace n connections by repeating the following n times:
 - (a) Randomly select an existing connection in the network.
 - (b) Randomly select a pair of neurons that are not connected. Do not select two neurons in the input layer or two neurons in the output layer.
 - (c) Remove the existing connection.
 - (d) Connect the pair of neurons that were not connected.

Note that this procedure is on average the same as the procedure described in section 2.3.1 with $p = 1 - \left(\frac{N-1}{N}\right)^n$, where N is the number of edges in the network; to contextualize the number of replaced connections, we will henceforth describe networks with our topology in terms of p and not n . We have also tried adding new connections without removing existing ones, and observed roughly the same performance; in this paper we replace connections to limit the number of parameters added to the network. We do not add connections within the input or output layers and we do not allow neurons to connect to themselves. We have seen good results with $p \approx 8\%$.

This topology allows us to train the network with a uniform

Network topology	Learning rate(s)	ϵ	β	N_{free}	$N_{weakly-clamped}$
Basic, unique learning rates	.128, .032, .008, .002	.5	1.0	500	8
Basic, one learning rate	.02	.5	1.0	500	8
Our topology, $p = 7.56\%$.02	.5	1.0	500	8

Table 1: Hyperparameters of networks tested on MNIST dataset

learning rate of .02 across the entire network. Weights of connections between pairs of neurons in adjacent layers are still initialized based on equation 10. Weights of intralayer connections and layer-skipping connections are drawn from the uniform distribution $U[-.05, .05]$, where the value .05 was chosen empirically to optimize network performance.

3.4 Tracking the training rates of individual pairs of layers

We find it informative to track the root-mean-square magnitude of corrections to weights in individual layers; if $\Delta w_{ij}^l(b)$ denotes the correction to the connection weight between the i^{th} neuron in layer l and the j^{th} neuron in layer $l + 1$ in response to batch b and N^l denotes the number of neurons in layer l , we define

$$\Delta w^l(b) := \sqrt{\frac{\sum_{i=1}^{N^l} \sum_{j=1}^{N^{l+1}} (\Delta w_{ij}^l(b))^2}{N^l N^{l+1}}} \quad (11)$$

as a metric of the extent to which the weights between layers l and $l + 1$ trained in response to batch b . We ignore the training of the weights of intralayer and layer-skipping connections.

This measurement is volatile, so for clarity we plot the average of $\Delta w^l(b')$ for b' in a neighborhood of batch b : for some n , we define

$$\Delta \hat{w}^l(b) := \frac{1}{2n+1} \sum_{b'=b-n}^{b+n} \Delta w^l(b'). \quad (12)$$

Comparing the traces $\Delta \hat{w}^l(b)$ vs. b for different layers reveals the extent of the vanishing gradient problem.

4 RESULTS

We compared the three network topologies described in sections 3.1, 3.2, and 3.3 on the MNIST dataset [11]. The hyperparameters for the networks are shown in table 1. All networks were trained for 250 epochs with 50,000 training examples, 10,000 test examples and a batch size of 20.

4.1 Network performance comparison

We tracked the error rates of the three networks as they were trained and found that the network with our topology significantly outperforms the basic network with one learning rate, and achieves the same error rates as the basic network with unique learning rates after around 25% more epochs. These results are shown in figure 2.

Notice that the basic network with unique learning rates converges to 0% training error and 2.5% test error in around 150 epochs and the network with our topology converges to 0% training error and 2.5% test error in around 190 epochs, whereas the basic network with one learning rate fails to converge in 250 epochs and

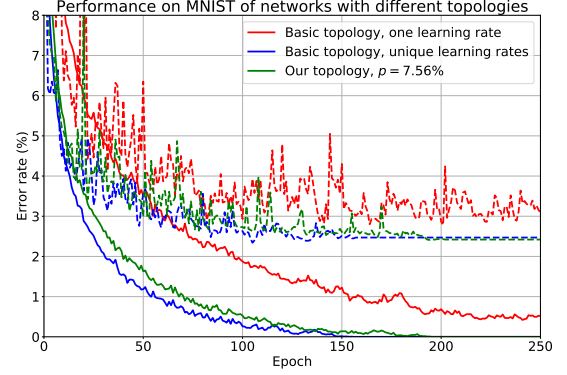


Figure 2: Comparison of performance of network topologies on MNIST dataset. Dotted lines show test error and solid lines show training error. In red is a network with the basic topology and a single learning rate (section 3.2). In blue is a network with the basic topology and unique learning rates (section 3.1), tuned to counter the vanishing gradient problem; this is a recreation of the 5-layer network in [17]. In green is a network with our topology, $p = 7.56\%$ (section 3.3). The network with a basic topology and a single learning rate performs poorly because it suffers from the vanishing gradient problem. The problem can be solved by introducing unique learning rates, or by implementing our topology.

has training and test error around .5% higher than the other two networks. While it is possible that the basic network with one learning rate would converge given enough time, it is clearly inferior to the other two networks. It is also apparent that in the context of the MNIST dataset, our network with a single learning rate is practically interchangeable with the basic network with unique learning rates.

4.2 Training rates of individual pairs of layers

We tracked the training rates of individual pairs of layers as described in equations 11 and 12 and found evidence that the extent of the vanishing gradient problem is the primary cause of the performance disparity seen in section 4.1. These results are shown in figure 3.

Notice in the basic network with one learning rate that there is a significant spread in the training rates of pairs of layers, with the deepest pair training at around 1% of the rate of the shallowest pair. This problem is solved very effectively in the basic network with unique learning rates. The problem appears to be solved effectively for the deepest 3 pairs in the network with our topology, but the output layer still trains significantly faster than the deeper 3 layers. This makes sense if we assume that the important factor in a layer's training rate is its expected path length to the target layer, because every neuron in the output layer connects to the target layer through a single connection, whereas paths starting at deeper neurons must first pass through the output layer before connecting to the target layer.

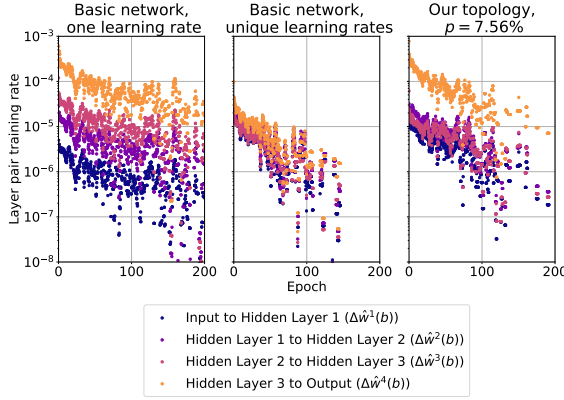


Figure 3: Extent of the vanishing gradient problem in different network topologies while training on MNIST dataset. To the left is a network with the basic topology and a single learning rate (section 3.2), in the center is a network with the basic topology and unique learning rates (section 3.1), and to the right is a network with our topology, 7.56% (section 3.3). In the basic topology with a single learning rate, the training rate of a layer decreases roughly by a factor of 4 each time its distance from the output increases by 1 layer. In the basic topology with unique learning rates, this effect is countered by increasing the learning rate of a layer by a factor of 4 each time its distance from the output increases by one layer; as a result, the training rates are very uniform. In the network with our topology, the effect is countered, albeit less-effectively, by replacing a subset of connections by random layer-skipping connections.

4.3 Error rate after one epoch as connections are added

We tracked the training error after one epoch of a network with our topology, with varying numbers of layer-skipping connections and found that the error rate decays approximately exponentially. These results are shown in figure 4.

The error rate drops quickly as connections are replaced early on. We believe this is because when a forward connection is added, in addition to providing a shortened path between the two neurons, it also provides paths via. at most 3 jumps to all pairs of neurons in the two layers, as a result of the intralayer connectivity. This benefit is exhausted when all pairs of layers are connected, leading to the more-gradual improvement later on. It appears to take around 25,000 replaced connections to reach this regime, likely because the extent of attenuation of a gradient depends less on the minimum path length between two neurons and more on the number of low-attenuation paths it has available. It appears that the error and spread do not asymptotically approach constant values as p approaches 1, which is unexpected; we do not currently have an explanation for this. There does not appear to be a close correspondence between a network’s spread and its error rate; while we believe that layer spread is a major contributor to the vanishing gradient problem, there are likely other contributors that we have

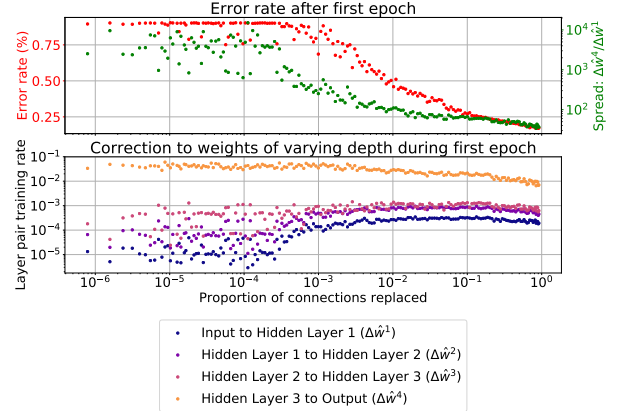


Figure 4: Performance of a network with our topology (section 3.3) with varying p . The top graph shows the training error after one epoch. The bottom graph shows the extent to which weights connecting each pair of layers was corrected over the epoch. It can be seen that there is little improvement for $p < 10^{-4}$, rapid improvement for $10^{-4} < p < 10^{-2}$ and moderate improvement for $p > 10^{-2}$. The training error after one epoch decreases as pairs of layers train more uniformly with respect to depth.

not recognized.

We found that our topology performs significantly worse than the basic topology with one learning rate when few connections are replaced, possibly due to inappropriate hyperparameter choices for intralayer connections. Some of our experiments have suggested that a network with our topology requires less care in initializing the weight matrix, but we have not explored this; if true, it could mean that improvement due to gradient flow is being conflated with improvement due to reduced impact of a poor weight initialization scheme.

5 DISCUSSION

5.1 Small-world metrics have little correlation with network performance

Contrary to our expectations, we found that a network’s performance does not track closely with its clustering coefficient or characteristic path length (appendix A), and that good performance is achieved when the algorithm in section 3.3 is executed with $p \approx 8\%$, which is substantially higher than the value of p that would make its topology small-world. We believe that our qualitative arguments about the merits of a small-world network are sound, but that the characteristic path length and clustering coefficient are poor predictors of the performance of a deep neural network. The characteristic path length does not take into account the number of low-attenuation paths a neuron can take to the output layer, but clearly this is an important factor, since a neuron’s behavior is dictated by an affine combination of all its inputs. While it is important for a network to retain its layered nature in order to generalize well, the clustering coefficient does not appear to be a good metric for this - e.g. a layered network without intralayer connections has a

small clustering coefficient because the neighbors of a given neuron are not connected to one-another.

There are several factors other than path length that we expect have an effect on performance as the value of p is changed, and it is possible that they disguise relationships with the small-world metrics. It has been shown [6] that in deep convolutional networks, performance can be improved by using layer-skipping connections to create a linear transformation (e.g. an identity mapping) in parallel with sets of nonlinear layers to reframe their task as learning the residual of the output of earlier layers; in our implementation layer-skipping connections form a linear transformation of earlier layers, and this fact may improve performance even though only some of the dimensions of the deeper layers are transferred. We discussed in section 2.2 that deep layers experience less evolution than shallow layers during the weakly-clamped phase; since layer-skipping connections allow deeper layers to begin to evolve sooner, we expect that they reduce this effect. As connections are replaced, a network looks less layered and more like a sparsely-connected single-layer network, so we expect that there is a tradeoff between the rate at which a network trains and its ability to generalize. In this paper our topology did not deteriorate test error relative to the original layered network, but it is highly possible that a network on a harder dataset (e.g. CIFAR or ImageNet) or with a higher value of p would see deterioration as its topology deviates from a layered topology.

5.2 Nonlinearities learning residuals

It has been shown [6, 9] that the performance of very-deep convolutional networks can be improved by using layer-skipping connections to create an identity mapping or other dimension-preserving linear transformation in parallel with groups of layers constituting nonlinear transformations. For a network without layer-skipping connections on a dataset $\{x_d, f(x_d)\}$, consider a group of layers. If the output of preceding layers is $g(x_d)$, the task of those layers is to learn $f(g^{-1}(x)) \approx x$ since the preceding layers have presumably trained to approximate f . Therefore, this task can be made easier by adding an identity mapping or linear transformation in parallel with the group of layers to reframe their task as finding the residual of the preceding layers $f(g^{-1}(x)) - x$, since it is easier to drive the output of a group of neurons towards zero than it is to make it approximate an arbitrary transformation.

Since layer-skipping connections constitute a linear transformation of a deep layer into a shallower layer, the above is a possible contributor to the effectiveness of layer-skipping connections; the MNIST dataset has a significant linear component, and of the hardened sigmoid activation function does a poor job of representing the underlying function, the network's effectiveness may be improved by skipping some of its hardened-sigmoid-based transformations using linear layer-skipping connections.

6 RELATED WORK

Much research has been done in pursuit of a biologically-plausible deep learning algorithm. [12], [25] and [16] are other algorithms that address the weight transport problem. [13] demonstrates that networks can train through backpropagation with randomly-initialized

feedback weights between neurons; the need for symmetric connection weights is not solved by equilibrium propagation and is not biologically-plausible, so this finding could potentially be applied here as well. [15] implements equilibrium propagation using spiking neurons; this type of activation function is believed to be more biologically-plausible than a hardened sigmoid function. [2] and [?] discuss the criteria a biologically-plausible algorithm would need to satisfy. [1] surveys promising biologically-motivated algorithms and evaluates their effectiveness on hard algorithms.

[19], [4] and [14] discuss neuromorphic architecture that could potentially implement equilibrium propagation as an analog computer.

Layer-skipping connections have been explored in other contexts. [6] and [22] use layer-skipping connections as a linear transformation in parallel with nonlinear layers to great effect in very-deep convolutional networks. [24] and [10] use a small-world topology in conventional multilayer feedforward networks.

The vanishing gradient problem was a big obstacle to the training of conventional deep networks through backpropagation, and [9] and [5] provide effective means for solving it in that context.

7 CONCLUSION

We believe that our topology is a suitable substitute for unique learning rates as a solution to the vanishing gradient problem. While it is not as effective, it is simpler, more biologically-plausible, and would be easier to implement in an analog computer.

There are several directions in which future research could be taken. It would be useful to find and mathematically-justify an effective weight initialization scheme for a network with our topology. It would be interesting to explore the effect of adding (instead of replacing) layer-skipping connections, training a network, then removing the connections and training it further, to see if they provide a residual benefit even after removal; doing so could allow a network in an analog computer to be trained quickly, and the added connections removed to reduce the power consumption and heat load of the network.

REFERENCES

- [1] Sergey Bartunov, Adam Santoro, Blake A. Richards, Geoffrey E. Hinton, and Timothy P. Lillicrap. Assessing the scalability of biologically-motivated deep learning algorithms and architectures. *CoRR*, abs/1807.04587, 2018.
- [2] Yoshua Bengio, Dong-Hyun Lee, Jörg Bornschein, and Zhouhan Lin. Towards biologically plausible deep learning. *CoRR*, abs/1502.04156, 2015.
- [3] E. Bullmore and O. Sporns. Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature*, 2009.
- [4] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Prasad Joshi, Andrew Lines, Andreas Wild, and Hong Wang. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, PP:1–1, 01 2018.
- [5] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In Yee Whye Teh and Mike Titterton, editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy, 13–15 May 2010. PMLR.
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. *CoRR*, abs/1512.03385, 2015.
- [7] John Hopfield. 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, 06 1984.
- [8] Mark Humphries and Kevin Gurney. Network ‘small-world-ness’: A quantitative method for determining canonical network equivalence. *PLoS one*, 3:e0002051, 02 2008.
- [9] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *CoRR*, abs/1502.03167, 2015.

- 2015.
- [10] Gokul Krishnan, Xiacong Du, and Yu Cao. Structural pruning in deep neural networks: A small-world approach, 2019.
 - [11] Y. LeCun and C. Cortes. The mnist database of handwritten digits, 1998.
 - [12] Dong-Hyun Lee, Saizheng Zhang, Asja Fischer, and Y. Bengio. Difference target propagation. pages 498–515, 08 2015.
 - [13] Timothy P. Lillicrap, Daniel Cownden, Douglas B. Tweed, and Colin J. Akerman. Random feedback weights support learning in deep neural networks, 2014.
 - [14] Mitchell Nahmias, Bhavin Shastri, A.N. Tait, and P.R. Prucnal. A leaky integrate-and-fire laser neuron for ultrafast cognitive computing. *Selected Topics in Quantum Electronics, IEEE Journal of*, 19:1–12, 09 2013.
 - [15] Peter O'Connor, Efstratios Gavves, and Max Welling. Training a network of spiking neurons with equilibrium propagation. 2018.
 - [16] Fernando J Pineda. Generalization of back-propagation to recurrent neural networks. *Physical Review Letters*, 59(19):2229–2232, 1987.
 - [17] Benjamin Scellier and Yoshua Bengio. Equilibrium propagation: Bridging the gap between energy-based models and backpropagation, 2016.
 - [18] Jürgen Schmidhuber. Deep learning in neural networks: An overview. *Neural Networks*, 61:85–117, Jan 2015.
 - [19] Jeffrey M. Shainline, Sonia M. Buckley, Adam N. McCaughan, Jeffrey T. Chiles, Amir Jafari Salim, Manuel Castellanos-Beltran, Christine A. Donnelly, Michael L. Schneider, Richard P. Mirin, and Sae Woo Nam. Superconducting optoelectronic loop neurons. *Journal of Applied Physics*, 126(4):044902, Jul 2019.
 - [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition, 2014.
 - [21] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Training very deep networks, 2015.
 - [22] Rupesh Kumar Srivastava, Klaus Greff, and Jürgen Schmidhuber. Highway networks. *CoRR*, abs/1505.00387, 2015.
 - [23] D. Watts and S. Strogatz. Collective dynamics of 'small-world' networks. *Nature*, 1998.
 - [24] L. Xiaohu, L. Xiaoling, Z. Jinhua, Z. Yulin, and L. Maolin. A new multilayer feedforward small-world neural network with its performances on function approximation. In *2011 IEEE International Conference on Computer Science and Automation Engineering*, 2011.
 - [25] Xiaohui Xie and HyunJune Seung. Equivalence of backpropagation and contrastive hebbian learning in a layered network. *Neural computation*, 15:441–54, 03 2003.

Appendix A METRICS OF A GRAPH'S SMALL-WORLDDNESS

The small-worldness of a graph is typically characterized by two metrics: a characteristic path length L and a clustering coefficient C [23]. A graph is small-world if L is small relative to the L of a random network with the same number of vertices and edges and C is not substantially smaller than the C of that random network; this situation can be quantified by a small-world coefficient σ [8]. The characteristic path length is the minimum number of edges in a path joining a pair of vertices, averaged over all pairs of vertices in the network. Specifically, for a graph with N vertices where $l(v_i, v_j)$ denotes the smallest number of edges needed to connect vertex v_i to vertex v_j ,

$$L = \frac{1}{N(N-1)} \sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N l(v_i, v_j). \quad (13)$$

The neighborhood of a vertex is the set of vertices with which it shares an edge. The clustering coefficient is the proportion of pairs of vertices in the neighborhood of a given vertex that share an edge, averaged over all vertices in the network. Specifically, for a graph with N vertices, if $N(v_i)$ denotes the neighborhood of vertex v_i and

$$c(v_j, v_k) = \begin{cases} 1 & v_k \in N(v_j) \\ 0 & \text{else} \end{cases},$$

$$C = \frac{1}{N} \sum_{i=1}^N \frac{1}{|N(v_i)|} \sum_{\substack{v_j \in N(v_i) \\ v_k \in N(v_i) \\ j \neq k}} c(v_j, v_k) \quad (14)$$

The small-worldness of a graph can be quantified by a small-world coefficient

$$\sigma = \frac{C/C_r}{L/L_r} \quad (15)$$

where C_r and L_r are the expected clustering coefficient and characteristic path length of a random graph with the same number of vertices and edges as the graph under consideration.

Appendix B INFLUENCE OF APPROXIMATION OF EQUATION OF MOTION ON VANISHING GRADIENT PROBLEM

One factor that probably contributes to the vanishing gradient problem is that when the output layer of neurons is perturbed at the beginning of the weakly-clamped phase, it takes an additional iteration of the approximation of equation 6 for the perturbation to influence the last hidden layer, then an additional iteration for the second-to-last hidden layer, and so on. We do not believe this factor to be a significant contributor to the vanishing gradient problem because increasing $N_{\text{weakly-clamped}}$ does not significantly affect network performance.

Assume that a typical neuron in the network changes by a small constant amount δ after each iteration of equation 6, and that the typical neuron is in state s_0 at the end of the free phase. After N iterations, a neuron in the output layer will have changed by $N\delta$ and one in a deep layer l layers from the output will have changed by $(N-l)\delta$. Assume further that s_0 is in the middle of the support of ρ' where ρ is the hardened sigmoid function from equation 1, so that we can replace $\rho(x)$ by x . Then it follows from equation 7 that a typical neuron in the output layer will be corrected by

$$\begin{aligned} \Delta w^{\text{output}} &= \frac{\alpha}{\beta} ((N\delta + s_0)((N-1)\delta + s_0) - s_0^2) \\ &\approx \frac{\alpha}{\beta} (N^2\delta^2 + 2s_0N\delta) \end{aligned} \quad (16)$$

and similarly a typical neuron from the deep layer will be corrected by

$$\Delta w^l \approx \frac{\alpha}{\beta} ((N-l)^2\delta^2 + 2s_0(N-l)\delta) \quad (17)$$

(where we have approximated $N \approx N-1$ and $N-l+1 \approx N-l$), so the ratio of these corrections will be

$$\begin{aligned} \frac{\Delta w^l}{\Delta w^{\text{output}}} &\approx \frac{(N-l)^2\delta^2 + 2s_0(N-l)\delta}{N^2\delta^2 + 2s_0N\delta} \\ &= \frac{N-l}{N} \frac{(N-l)\delta + 2s_0}{N\delta + 2s_0} \\ &\approx 1 - \frac{l}{N} \end{aligned} \quad (18)$$

(where we have approximated $\frac{(N-l)\delta+2s_0}{N\delta+2s_0} \approx 1$). We have observed that

$$\frac{\Delta w^l}{\Delta w^{\text{output}}} \sim 4^{-l}, \quad (19)$$

for the type of networks we tested, so it seems unlikely that this phenomenon is a primary cause of the vanishing gradient problem.