

---

# Neural Synaptic Balance

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 For a given additive cost function  $R$  (regularizer), a neuron is said to be in balance  
2 if the total cost of its input weights is equal to the total cost of its output weights.  
3 The basic example is provided by feedforward layered networks of ReLU units  
4 trained with  $L_2$  regularizers, which exhibit balance after proper training. We  
5 develop a general theory that extends this phenomenon in three broad directions  
6 in terms of: (1) activation functions; (2) regularizers, including all  $L_p$  ( $p > 0$ )  
7 regularizers; and (3) architectures (non-layered, recurrent, convolutional, mixed  
8 activations). Gradient descent on the error function alone does not converge in  
9 general to a balanced state where every neuron is in balance, even when starting  
10 from a balanced state. However, gradient descent on the regularized error function  
11 must converge to a balanced state, and thus network balance can be used to assess  
12 learning progress. The theory is based on two local neuronal operations: scaling  
13 which is commutative, and balancing which is not commutative. Finally, and most  
14 importantly, given any initial set of weights, when local balancing operations are  
15 applied to each neuron in a stochastic manner, global order always emerges through  
16 the convergence of the stochastic algorithm to the same unique set of balanced  
17 weights. The reason for this convergence is the existence of an underlying strictly  
18 convex optimization problem where the relevant variables are constrained to a  
19 linear, only architecture-dependent, manifold. The theory is corroborated through  
20 simulations carried out on benchmark data sets. Balancing operations are entirely  
21 local and thus physically plausible in biological and neuromorphic networks.

## 22 1 Introduction

23 When large neural networks are trained on complex tasks, they produce large arrays of synaptic  
24 weights that have no clear structure and are difficult to interpret. Thus finding any kind of structure in  
25 the weights of large neural networks is of great interest. Here we study a particular kind of structure  
26 we call neural synaptic balance and the conditions under which it emerges. Neural synaptic balance  
27 is different from the biological notion of balance between excitation and inhibition [Froemke, 2015,  
28 Field et al., 2020, Howes and Shatalina, 2022, Kim and Lee, 2022, Shirani and Choi, 2023]. We  
29 use this term to refer to any systematic relationship between the input and output synaptic weights  
30 of individual neurons or layers of neurons. Here we consider the case where the cost of the input  
31 weights is equal to the cost of the output weights, where the cost is defined by some regularizer. One  
32 of the most basic examples of such a relationship is when the sum of the squares of the input weights  
33 of a neuron is equal to the sum of the squares of its output weights.

34 **Basic Example:** The basic example where this happens is with a neuron with a ReLU activation  
35 function inside a network trained to minimize an error function with  $L_2$  regularization. If we multiply  
36 the incoming weights of the neuron by some  $\lambda > 0$  (including the bias) and divide the outgoing  
37 weights of the neuron by the same  $\lambda$ , it is easy to see that this scaling operation does not affect in any  
38 way the contribution of the neuron to the rest of the network. Thus, the component of the overall

39 error function that depends only on the input-output function of the network is unchanged. However,  
 40 the value of the  $L_2$  regularizer changes with  $\lambda$  and we can ask what is the value of  $\lambda$  that minimizes  
 41 the corresponding contribution given by:

$$\sum_{i \in IN} (\lambda w_i)^2 + \sum_{i \in OUT} (w_i/\lambda)^2 = \lambda^2 A + \frac{1}{\lambda^2} B \quad (1.1)$$

42 where  $IN$  and  $OUT$  denote the set of incoming and outgoing weights respectively,  $A = \sum_{i \in IN} w_i^2$ ,  
 43 and  $B = \sum_{i \in OUT} w_i^2$ . The product of the two terms on the right-hand side of Equation 1.1 is equal to  
 44  $AB$  and does not depend on  $\lambda$ . Thus, the minimum is achieved when these two terms are equal, which  
 45 yields:  $(\lambda^*)^4 = B/A$  for the optimal  $\lambda^*$ . The corresponding new set of weights,  $v_i = \lambda^* w_i$  for the  
 46 input weights and  $v_i = w_i/\lambda^*$  for the outgoing weights, must be balanced:  $\sum_{i \in IN} v_i^2 = \sum_{i \in OUT} v_i^2$ .  
 47 This is because its optimal scaling factor can only be  $\lambda^* = 1$ . Thus, we can define two operations  
 48 that can be applied to the incoming and outgoing weights of a neuron: scaling and balancing. It  
 49 is easy to check that scaling operations applied to any two neurons commute, whereas balancing  
 50 operations do not commute if the two neurons are directly connected (Appendix). If a network of  
 51 ReLU neurons is properly trained using a standard error function with an  $L_2$  regularizer, at the end of  
 52 training one observes a remarkable phenomenon: for each ReLU neuron, the norm of the incoming  
 53 synaptic weights is approximately equal to the norm of the outgoing synaptic weights, i.e. every  
 54 neuron is balanced.

55 There have been isolated previous studies of this kind of synaptic balance [Du et al., 2018, Stock  
 56 et al., 2022] under special conditions. For instance, in Du et al. [2018], it is shown that if a deep  
 57 network is initialized in a balanced state with respect to the sum of squares metric, and if training  
 58 progresses with an infinitesimal learning rate, then balance is preserved throughout training. Here,  
 59 we take a different approach aimed at uncovering the generality of neuronal balance phenomena,  
 60 the learning conditions under which they occur, as well as new local balancing algorithms and their  
 61 convergence properties. We study neural synaptic balance in its generality in terms of activation  
 62 functions, regularizers, network architectures, and training stages. In particular, we systematically  
 63 answer questions such as: Why does balance occur? Does it occur only with ReLU neurons? Does it  
 64 occur only with  $L_2$  regularizers? Does it occur only in fully connected feedforward architectures?  
 65 Does it occur only at the end of training? And what happens if we balance neurons at random in a  
 66 large network?

## 67 2 Generalization of the Activation Functions

68 What enables scaling ReLU neurons without changing their input-output function is the homogeneous  
 69 property of ReLU activation function. An activation function  $f$  is said to be *homogeneous* if for every  
 70  $\lambda > 0$ ,  $f(\lambda x) = \lambda f(x)$ . To fully characterize the class of homogeneous activation functions, we first  
 71 define a new class of activation functions, corresponding to bilinear units (BiLU), consisting of two  
 72 half-lines meeting at the origin.

73 **Definition 2.1.** (BiLU) A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is bilinear (BiLU) if and only if  
 74  $f(x) = ax$  when  $x < 0$ , and  $f(x) = bx$  when  $x \geq 0$ , for some fixed parameters  $a$  and  $b$  in  $\mathbb{R}$ .

75 BiLU units include linear units ( $a = b$ ), ReLU units ( $a = 0, b = 1$ ), leaky ReLU ( $a = \epsilon; b = 1$ ) units,  
 76 and symmetric linear units ( $a = -b$ ), all of which can also be viewed as special cases of piece-wise  
 77 linear units [Tavakoli et al., 2021], with a single hinge. One advantage of ReLU and more generally  
 78 BiLU neurons, which is very important during backpropagation learning, is that their derivative is  
 79 very simple and can only take one of two values ( $a$  or  $b$ ). We have the following equivalence.

80 **Proposition 2.2.** A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is homogeneous if and only if it is a  
 81 BiLU activation function.

82 *Proof.* Every function in BiLU is clearly homogeneous. Conversely, any homogeneous function  $f$   
 83 must satisfy: (1)  $f(0x) = 0f(x) = f(0) = 0$ ; (2)  $f(x) = f(1x) = f(1)x$  for any positive  $x$ ; and (3)  
 84  $f(x) = f(-u) = f(-1)u = -f(-1)x$  for any negative  $x$ . Thus  $f$  is in BiLU with  $a = -f(-1)$   
 85 and  $b = f(1)$ .  $\square$

86 In the Appendix, we provide a simple proof that networks of BiLU neurons, even with a single hidden  
 87 layer, have universal approximation properties.

88 While in the rest of this work we use BiLU neurons, it is possible to generalize the notions of scaling  
 89 and balancing even further. To see this, suppose that there is a neuron with an activation function  
 90  $f : \mathbb{R} \rightarrow \mathbb{R}$ , and functions  $g : (a, b) \rightarrow \mathbb{R}$  and  $h : (a, b) \rightarrow \mathbb{R}$ , such that:  $f(g(\lambda)x) = h(\lambda)f(x)$ ,  
 91 for any  $\lambda \in (a, b)$ . Then if we multiply the incoming weights by  $g(\lambda)$  and divide the outgoing  
 92 weights by  $h(\lambda) \neq 0$  (generalized scaling), we see again that the influence of the neuron on the  
 93 rest of the network is unchanged. And thus, again, we can try to find the value of  $\lambda$  that minimizes  
 94 the regularization cost (generalized balancing). Here we provide an example of such an activation  
 95 function, with  $g(\lambda) = \lambda$  and  $h(\lambda) = \lambda^c$ . Additional details are given in the Appendix.

96 **Proposition 2.3.** *The set of activation functions  $f$  satisfying  $f(\lambda x) = \lambda^c f(x)$  for any  $x \in \mathbb{R}$  and  
 97 any  $\lambda > 0$  consist of the functions of the form:*

$$f(x) = \begin{cases} Cx^c & \text{if } x \geq 0 \\ Dx^c & \text{if } x < 0. \end{cases} \quad (2.1)$$

98 where  $c \in \mathbb{R}$ ,  $C = f(1) \in \mathbb{R}$ , and  $D = f(-1) \in \mathbb{R}$ . We call these bi-power units (BiPU). If, in  
 99 addition, we want  $f$  to be continuous at 0, we must have either  $c > 0$ , or  $c = 0$  with  $C = D$ .

100 Note that in the general case where  $c > 0$ ,  $C$  and  $D$  do not need to be equal. In particular, one of  
 101 them can be equal to zero, and the other one can be different from zero giving rise to rectified power  
 102 units.

### 103 3 Generalization of the Regularizers

104 As we have seen, given a BiLU neuron, scaling its input and output weights by  $\lambda$  and  $1/\lambda$  respectively  
 105 does not alter its contribution to the rest of the network and thus we can adjust  $\lambda$  to reduce or even  
 106 minimize the contribution of the corresponding weights to the regularizer. It is reasonable to assume  
 107 that the regularizer has the general additive form:  $R(W) = \sum_w g_w(w)$  where  $W$  denotes all the  
 108 weights in the network. Without much loss of generality, we can assume that the  $g_w$  are continuous,  
 109 and lower-bounded by 0. To ensure the existence and uniqueness of a minimum during the balancing  
 110 of any neuron, We will assume that each function  $g_w$  depends only on the magnitude  $|w|$  of the  
 111 corresponding weight, and that  $g_w$  monotonically increases from 0 to  $+\infty$ . Clearly,  $L_2, L_1$  and  
 112 more generally all  $L_p$  regularizers are special cases where, for  $p > 0$ ,  $L^p$  regularization is defined  
 113 by:  $R(W) = \sum_w |w|^p$ . Differentiability conditions can be added to be able to derive closed form  
 114 solutions for the balance (optimal scaling). This is satisfied by all forms of  $L_p$  regularization, for  
 115  $p > 0$ . We have the following theorem.

116 **Theorem 3.1.** *(Balance and Regularizer Minimization) Assume an additive regularizer with the  
 117 properties described above, where in addition we assume that the functions  $g_w$  are continuously  
 118 differentiable, except perhaps at the origin. Then, for any neuron, there exists one optimal value  $\lambda^*$   
 119 that minimizes  $R(W)$ . This value must be a solution of the consistency equation:*

$$\lambda^2 \sum_{w \in IN(i)} w g'_w(\lambda w) = \sum_{w \in OUT(i)} w g'_w(w/\lambda) \quad (3.1)$$

120 *Once the weights are rebalanced accordingly, the new weights must satisfy the generalized balance  
 121 equation:*

$$\sum_{w \in IN(i)} w g'(w) = \sum_{w \in OUT(i)} w g'(w) \quad (3.2)$$

122 *In particular, if  $g_w(w) = |w|^p$  for all the incoming and outgoing weights of neuron  $i$ , then the optimal  
 123 value  $\lambda^*$  is unique and equal to:*

$$\lambda^* = \left( \frac{\sum_{w \in OUT(i)} |w|^p}{\sum_{w \in IN(i)} |w|^p} \right)^{1/2p} = \left( \frac{\|OUT(i)\|_p}{\|IN(i)\|_p} \right)^{1/2} \quad (3.3)$$

124 *After balancing, the decrease  $\Delta R \geq 0$  in the value of the  $L_p$  regularizer  $R = \sum_w |w|^p$  is given by:*

$$\Delta R = \left( \left( \sum_{w \in IN(i)} |w|^p \right)^{1/2} - \left( \sum_{w \in OUT(i)} |w|^p \right)^{1/2} \right)^2 \quad (3.4)$$

125 After balancing neuron  $i$ , its new weights satisfy the generalized  $L_p$  balance equation:

$$\sum_{w \in IN(i)} |w|^p = \sum_{w \in OUT(i)} |w|^p \quad (3.5)$$

126 *Proof.* The results are obtained by setting the derivative of the regularizer with respect to the scaling  
 127 factor  $\lambda$  to 0. Note that the theorem applies to regularizers combining different  $L_p$ 's (e.g. of the form  
 128  $\alpha L_2 + \beta L_1$ ). The details are given in the Appendix.  $\square$

## 129 4 Generalization of the Architectures

130 It is straightforward to check that the scaling and balancing operations can be extended in the  
 131 following cases (see Appendix for additional details):

- 132 1. Mixed networks containing both BiLU and non-BiLU units. One can just restrict those  
 133 operations to the BiLU neurons.
- 134 2. Recurrent networks containing BiLU neurons, not just feedforward networks.
- 135 3. Networks that are not layered, or not fully connected.
- 136 4. In addition, scaling and balancing operations can be applied layer-wise to an entire layer of  
 137 BiLU neurons in a tied manner, by using the same scaling factor  $\lambda$  with a single optimal  
 138 value  $\lambda^*$  for all the neurons in the layer. In particular, this allows the application of scaling  
 139 and balancing to convolutional layers of BiLU neurons.

## 140 5 Balancing Algorithms

141 **Gradient Descent:** When a network of BiLU neurons is trained by gradient descent to minimize  
 142 an error function  $E(W)$ , such as the negative log-likelihood of the data, there is no reason for the  
 143 final weights to be balanced. However, when a network is properly trained to minimize a regularized  
 144 error function  $\mathcal{E} = E(W) + R(W)$ , the final weights ought to be balanced. The reason is that if a  
 145 neuron is not in a balanced state at the end of training, then we can further reduce its contribution to  
 146  $R$  smoothly by balancing it. This implies that the gradient of  $\mathcal{E}(W)$  is not equal to zero at the end of  
 147 training, and thus training has not properly converged. The converse is that the degree of balance can  
 148 be used as a proxy for assessing whether learning has converged or not.

149 **Stochastic Balancing:** More interestingly, we now investigate what happens if we fix the weights  $W$   
 150 of a network and iteratively balance its BiLU neurons.

151 **Theorem 5.1.** (*Convergence of Stochastic Balancing*) Consider a network of BiLU neurons with  
 152 an error function  $\mathcal{E}(W) = E(W) + R(W)$  where  $R$  is any  $L_p$  ( $p > 0$ ) regularizer. Let  $W$  denote  
 153 the initial weights. When the neuronal stochastic balancing algorithm is applied throughout the  
 154 network so that every neuron is visited from time to time, then  $E(W)$  remains unchanged but  $R(W)$   
 155 must converge to some finite value that is less or equal to the initial value, strictly less if the initial  
 156 weights are not balanced. In addition, for every neuron  $i$ ,  $\lambda_i^*(t) \rightarrow 1$  and the weights themselves must  
 157 converge to a limit  $W^*$  which is globally balanced, with  $E(W) = E(W^*)$  and  $R(W) \geq R(W^*)$ ,  
 158 and with equality if only if  $W$  is already balanced. Finally,  $W^*$  is unique as it corresponds to the  
 159 solution of a strictly convex optimization problem with special linear constraints that depend only on  
 160 the network architecture (and not on  $W$ ). Stochastic balancing projects to stochastic trajectories in  
 161 the linear manifold that run from the origin to the unique optimal configuration.

162 *Proof.* Each individual balancing operation leaves  $E(W)$  unchanged because the BiLU neurons are  
 163 homogeneous. Furthermore, each balancing operation reduces the regularization error  $R(W)$ , or  
 164 leaves it unchanged. Since the regularizer is lower-bounded by zero, the value of the regularizer must  
 165 approach a limit as the stochastic updates are being applied. However, this alone does not imply

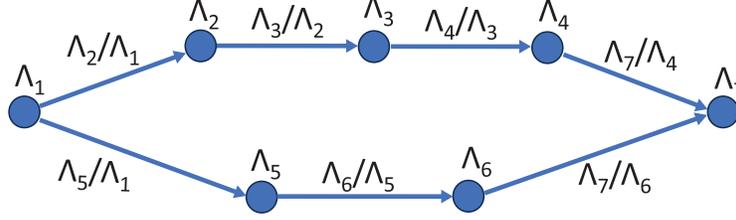


Figure 1: Two hidden units (1 and 7) connected by two different directed paths 1-2-3-4-7 and 1-5-6-7 in a BiLU network. Each unit  $i$  has a scaling factor  $\Lambda_i$ , and each directed edge from unit  $j$  to unit  $i$  has a scaling factor  $M_{ij} = \Lambda_i/\Lambda_j$ . The products of the  $M_{ij}$ 's along each path is equal to:  $\frac{\Lambda_2}{\Lambda_1} \frac{\Lambda_3}{\Lambda_2} \frac{\Lambda_4}{\Lambda_3} \frac{\Lambda_7}{\Lambda_4} = \frac{\Lambda_5}{\Lambda_1} \frac{\Lambda_6}{\Lambda_5} \frac{\Lambda_7}{\Lambda_6} = \frac{\Lambda_7}{\Lambda_1}$ . Therefore the variables  $L_{ij} = \log M_{ij}$  must satisfy the linear equation:  $L_{21} + L_{32} + L_{43} + L_{74} = L_{51} + L_{65} + L_{76} = \log \Lambda_7 - \log \Lambda_1$ .

166 that the weights are converging and whether the limit is unique or not. To address these issues, for  
 167 simplicity, we use a continuous time notation. After a certain time  $t$  each neuron has been balanced a  
 168 certain number of times. While the balancing operations are not commutative as balancing operations,  
 169 they are commutative as scaling operations. Thus we can reorder the scaling operations and group  
 170 them neuron by neuron so that, for instance, neuron  $i$  has been scaled by the sequence of scaling  
 171 operations of the form:

$$S_{\lambda_1^*}(i) S_{\lambda_2^*}(i) \dots S_{\lambda_{n_{it}}^*}(i) = S_{\Lambda_i(t)}(i) \quad (5.1)$$

172 where  $n_{it}$  corresponds to the count of the last update of neuron  $i$  prior to time  $t$ , and:

$$\Lambda_i(t) = \prod_{1 \leq n \leq n_{it}} \lambda_n^*(i) \quad (5.2)$$

173 For the input and output units, we can consider that their balancing coefficients  $\lambda^*$  are always equal  
 174 to 1 (at all times) and therefore  $\Lambda_i(t) = 1$  for any visible unit  $i$ . At time  $t$  the weight connecting unit  
 175  $j$  to unit  $i$  is given by:  $w_{ij}(t) = w_{ij}(0) \Lambda_i(t) / \Lambda_j(t)$ , where  $w_{ij}(0)$  corresponds to the initial value.  
 176 In the Appendix, we show upfront that for all BiLU units  $i$ ,  $\Lambda_i(t)$  converges to some limit  $\Lambda_i > 0$ ,  
 177 and thus the weights converge too. Here, we first suppose that the coefficients  $\Lambda_i(t)$  converge to  
 178 some limit  $\Lambda_i$ , and recover the convergence at the end from understanding the overall proof. As a  
 179 result, for any  $L_p$  regularizer, the coefficients  $\Lambda_i$  corresponding to a globally balanced state must be  
 180 solutions of the following optimization problem:

$$\min_{\Lambda} R(\Lambda) = \sum_{ij} \left| \frac{\Lambda_i}{\Lambda_j} w_{ij} \right|^p \quad (5.3)$$

181 under the simple constraints:  $\Lambda_i > 0$  for all the BiLU hidden units, and  $\Lambda_i = 1$  for all the visible (input  
 182 and output) units. In this form, the problem is not convex. Introducing new variables  $M_j = 1/\Lambda_j$   
 183 is not sufficient to render the problem convex. Using variables  $M_{ij} = \Lambda_i/\Lambda_j$  is better, but still  
 184 problematic for  $0 < p \leq 1$ . However, let us instead introduce the new variables  $L_{ij} = \log(\Lambda_i/\Lambda_j)$ .  
 185 These are well defined since we know that  $\Lambda_i/\Lambda_j > 0$ . The objective now becomes:

$$\min R(L) = \sum_{ij} |e^{L_{ij}} w_{ij}|^p = \sum_{ij} e^{pL_{ij}} |w_{ij}|^p \quad (5.4)$$

186 This objective is strictly convex in the variables  $L_{ij}$ , as a sum of strictly convex functions (exponen-  
 187 tials). However, to show that it is a convex optimization problem we need to study the constraints  
 188 on the variables  $L_{ij}$ . In particular, from the set of  $\Lambda_i$ 's it is easy to construct a unique set of  $L_{ij}$ .  
 189 However what about the converse?

190 **Definition 5.2.** A set of real numbers  $L_{ij}$ , one per connection of a given neural architecture, is  
 191 self-consistent if and only if there is a unique corresponding set of numbers  $\Lambda_i > 0$  (one per unit)  
 192 such that:  $\Lambda_i = 1$  for all visible units and  $L_{ij} = \log \Lambda_i/\Lambda_j$  for every directed connection from a unit  
 193  $j$  to a unit  $i$ .

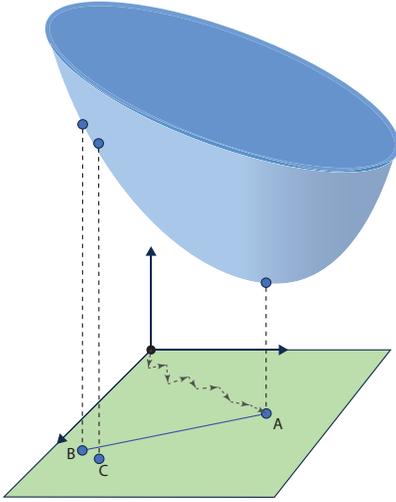


Figure 2: The problem of minimizing the strictly convex regularizer  $R(L_{ij}) = \sum_{ij} e^{pL_{ij}} |w_{ij}|^p$  ( $p > 0$ ), over the linear (hence convex) manifold of self-consistent configurations defined by the linear constraints of the form  $\sum_{\pi} L_{ij} = 0$ , where  $\pi$  runs over input-output paths. The regularizer function depends on the weights. The linear manifold depends only on the architecture, i.e., the graph of connections. This is a strictly convex optimization problem with a unique solution associated with the point  $A$ . At  $A$  the corresponding weights must be balanced, or else a self-consistent configuration of lower cost could be found by balancing any non-balanced neuron. Finally, any other self-consistent configuration  $B$  cannot correspond to a balanced state of the network, since there must exist balancing moves that further reduce the regularizer cost (see main text). Stochastic balancing produces random paths from the origin, where  $L_{ij} = \log M_{ij} = 0$ , to the unique optimum point  $A$ .

194 *Remark 5.3.* This definition depends on the graph of connections, but not on the original values of  
 195 the synaptic weights. Every balanced state is associated with a self-consistent set of  $L_{ij}$ , but not  
 196 every self-consistent set of  $L_{ij}$  is associated with a balanced state.

197 **Proposition 5.4.** *A set  $L_{ij}$  associated with a neural architecture is self-consistent if and only if*  
 198  *$\sum_{\pi} L_{ij} = 0$  where  $\pi$  is any directed path connecting an input unit to an output unit or any directed*  
 199 *cycle (for recurrent networks).*

200 *Proof.* If we look at any directed path  $\pi$  from unit  $i$  to unit  $j$ , it is easy to see that we must have:

$$\sum_{\pi} L_{kl} = \log \Lambda_i - \log \Lambda_j \quad (5.5)$$

201 This is illustrated in Figure 1. Thus along any directed path that connects any input unit to any output  
 202 unit, we must have  $\sum_{\pi} L_{ij} = 0$ . In addition, for recurrent neural networks, if  $\pi$  is a directed cycle  
 203 we must also have:  $\sum_{\pi} L_{ij} = 0$ . Thus in short we only need to add linear constraints of the form:  
 204  $\sum_{\pi} L_{ij} = 0$ . Any unit is situated on a path from an input unit to an output unit. Along that path, it is  
 205 easy to assign a value  $\Lambda_i$  to each unit by simple propagation starting from the input unit which has a  
 206 multiplier equal to 1. When the propagation terminates in the output unit, it terminates consistently  
 207 because the output unit has a multiplier equal to 1 and, by assumption, the sum of the multipliers  
 208 along the path must be zero. So we can derive scaling values  $\Lambda_i$  from the variables  $L_{ij}$ . Finally, it is  
 209 easy to show that there are no clashes, i.e. that it is not possible for two different propagation paths to  
 210 assign different multiplier values to the same unit  $i$  (see Appendix).  $\square$

211 *Remark 5.5.* Thus the constraints associated with being a self-consistent configuration of  $L_{ij}$ 's are  
 212 all linear. This linear manifold of constraints depends only on the architecture, i.e., the graph of  
 213 connections. The strictly convex function  $R(L_{ij})$  depends on the actual weights  $W$ . Different sets of  
 214 weights  $W$  produce different convex functions over the same linear manifold.

215 *Remark 5.6.* One could coalesce all the input units and all output units into a single unit, in which  
 216 case a path from an input unit to and output unit becomes also a directed cycle. In this representation,  
 217 the constraints are that the sum of the  $L_{ij}$  must be zero along any directed cycle. In general, it is not  
 218 necessary to write a constraint for every path from input units to output units. It is sufficient to select  
 219 a representative set of paths such that every unit appears in at least one path.

220 We can now complete the proof of Theorem 5.1. Given a neural network of BiLUs with a set  
 221 of weights  $W$ , we can consider the problem of minimizing the regularizer  $R(L_{ij})$  over the self-  
 222 admissible configuration  $L_{ij}$ . For any  $p > 0$ , the  $L_p$  regularizer is strictly convex and the space of  
 223 self-admissible configurations is linear and hence convex. Thus this is a strictly convex optimization

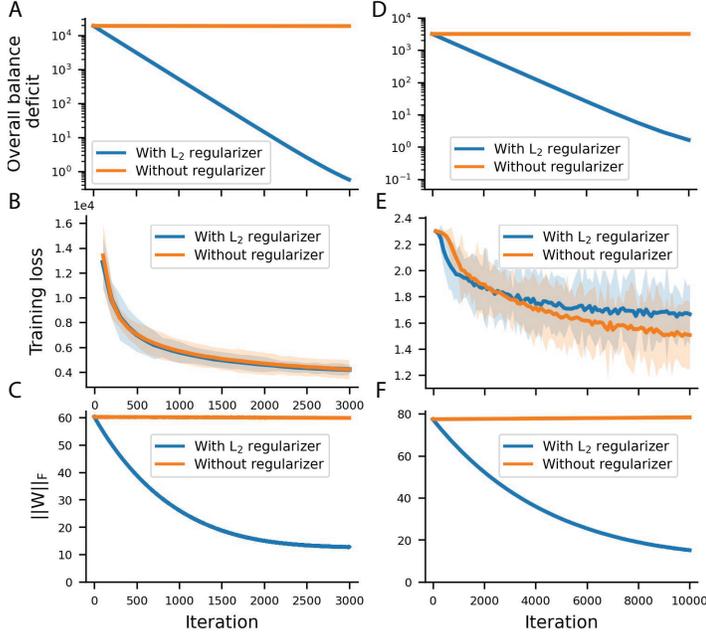


Figure 3: **SGD applied to  $E$  alone, in general, does not converge to a balanced state, but SGD applied to  $E + R$  converges to a balanced state.** (A-C) Simulations use a deep fully connected autoencoder trained on the MNIST dataset. (D-F) Simulations use a deep locally connected network trained on the CFAR10 dataset. (A,D) Regularization leads to neural balance. (B,E) The training loss decreases and converges during training (these panels are not meant for assessing the quality of learning when using a regularizer). (C,F) Using weight regularization decreases the norm of weights. (A-F) Shaded areas correspond to one s.t.d around the mean (in some cases the s.t.d. is small and the shaded area is not visible).

224 problem that has a unique solution (Figure 2). Note that the minimization is carried over self-  
 225 consistent configurations, which in general are not associated with balanced states. However, the  
 226 configuration of the weights associated with the optimum set of  $L_{ij}$  (point  $A$  in Figure 2) must be  
 227 balanced. To see this, imagine that one of the BiLU units—unit  $i$  in the network is not balanced. Then  
 228 we can balance it using a multiplier  $\lambda_i^*$  and replace  $\Lambda_i$  by  $\Lambda_i' = \Lambda_i \lambda_i^*$ . It is easy to check that the new  
 229 configuration including  $\Lambda_i'$  is self-consistent. Thus, by balancing unit  $i$ , we are able to reach a new  
 230 self-consistent configuration with a lower value of  $R$  which contradicts the fact that we are at the  
 231 global minimum of the strictly convex optimization problem.

232 We know that the stochastic balancing algorithm always converges to a balanced state. We need to  
 233 show that it cannot converge to any other balanced state, and in fact that the global optimum is the  
 234 only balanced state. By contradiction, suppose it converges to a different balanced state associated  
 235 with the coordinates  $(L_{ij}^B)$  (point  $B$  in Figure 2). Because of the self-consistency, this point is also  
 236 associated with a unique set of  $(\Lambda_i^B)$  coordinates. The cost function is continuous and differentiable  
 237 in both the  $L_{ij}$ 's and the  $\Lambda_i$ 's coordinates. If we look at the negative gradient of the regularizer, it  
 238 is non-zero and therefore it must have at least one non-zero component  $\partial R / \partial \Lambda_i$  along one of the  
 239  $\Lambda_i$  coordinates. This implies that by scaling the corresponding unit  $i$  in the network, the regularizer  
 240 can be further reduced, and by balancing unit  $i$  the balancing algorithm will reach a new point ( $C$  in  
 241 Figure 2) with lower regularizer cost. This contradicts the assumption that  $B$  was associated with a  
 242 balanced state. Thus, given an initial set of weights  $W$ , the stochastic balancing algorithm must  
 243 always converge to the same and unique optimal balanced state  $W^*$  associated with the self-consistent  
 244 point  $A$ . A particular stochastic schedule corresponds to a random path within the linear manifold  
 245 from the origin (at time zero, all the multipliers are equal to 1, and therefore  $M_{ij} = 1$  and  $L_{ij} = 0$   
 246 for any  $i$  and any  $j$ ) to the unique optimum point  $A$ .

247

□

248 *Remark 5.7.* From the proof, it is clear that the same result holds also for any deterministic balancing  
 249 schedule, as well as for tied and non-tied subset balancing, e.g., for layer-wise balancing and tied  
 250 layer-wise balancing. In the Appendix, we provide an analytical solution for the case of tied layer-wise  
 251 balancing in a layered feed-forward network.

252 *Remark 5.8.* From the proof, it is also clear that the same convergence to the unique global optimum  
 253 is observed if each neuron, when stochastically visited, is favorably scaled rather than balanced, i.e.,  
 254 it is scaled with a factor that reduces  $R$  but not necessarily minimizes  $R$ . Stochastic balancing can  
 255 also be viewed as a form of EM algorithm where the E and M steps can be taken fully or partially.

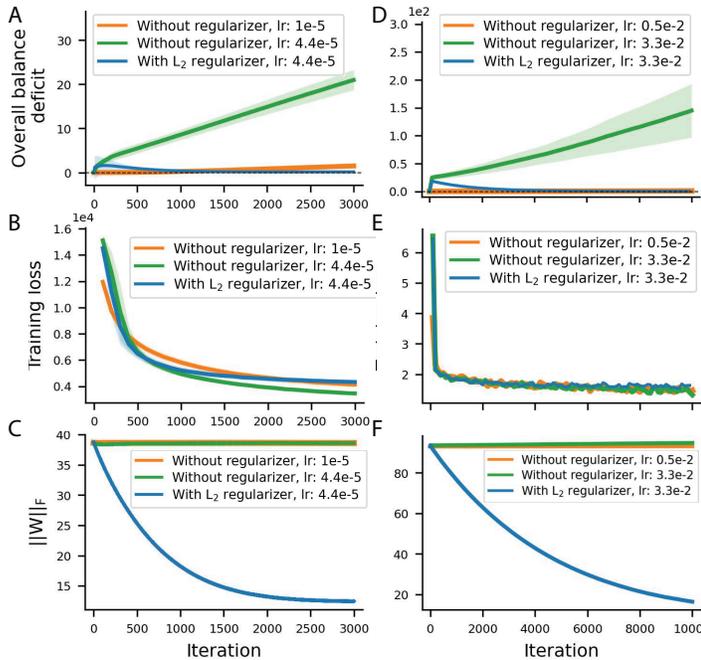


Figure 4: **Even if the starting state is balanced, SGD does not preserve the balance unless the learning rate is infinitely small.** (A-C) Simulations use a deep fully connected autoencoder trained on the MNIST dataset. (D-F) Simulations use a deep locally connected network trained on the CFAR10 dataset. (A-F) The initial weights are balanced using the stochastic balancing algorithm. Then the network is trained by SGD. (A,D) When the learning rate (lr) is relatively large, without regularization, the initial balance of the network is rapidly disrupted. (B,E) The training loss decreases and converges during training (these panels are not meant for assessing the quality of learning when using a regularizer). (C,F) Using weight regularization decreases the norm of the weights. (A-F) Shaded areas correspond to one s.t.d. around the mean (in some cases the s.t.d. is small and the shaded area is not visible).

## 256 6 Simulations

257 To further corroborate the results, we ran multiple experiments. Here we report the results from two  
 258 series of experiments. The first one is conducted using a six-layer, fully connected, autoencoder  
 259 trained on MNIST [Deng, 2012] for a reconstruction task with ReLU activation functions in all layers  
 260 and the sum of squares errors loss function. The number of neurons in consecutive layers, from  
 261 input to output, is 784, 200, 100, 50, 100, 200, 784. Stochastic gradient descent (SGD) learning by  
 262 backpropagation is used for learning with a batch size of 200.

263 The second one is conducted using three locally connected layers followed by three fully connected  
 264 layers trained on CFAR10 [Krizhevsky and Hinton, 2009] for a classification task with leaky ReLU  
 265 activation functions in the hidden layers, a softmax output layer, and the cross entropy loss function.  
 266 The number of neurons in consecutive layers, from input to output, is 3072, 5000, 2592, 1296, 300,  
 267 100, 10. Stochastic gradient descent (SGD) learning by backpropagation is used for learning with a  
 268 batch size of 5.

269 In all the simulation figures (Figures 3, 4, and 5) the left column presents results obtained from the  
 270 first experiment, while the right column presents results obtained from the second experiment. While  
 271 we used both  $L_1$  and  $L_2$  regularizers in the experiments, in the figures we report the results obtained  
 272 with the  $L_2$  regularizer, which is the most widely used regularizer. In Figures 3 and 4, training is  
 273 done using batch gradient descent on the MNIST and CIFAR data. The balance deficit for a single  
 274 neuron  $i$  is defined as:  $(\sum_{w \in IN(i)} w^2 - \sum_{w \in OUT(i)} w^2)^2$ , and the overall balance deficit is defined  
 275 as the sum of these single-neuron balance deficits across all the hidden neurons in the network. The  
 276 overall deficit is zero if and only if each neuron is in balance. In all the figures,  $\|W\|_F$  denotes the  
 277 Frobenius norm of the weights.

278 Figure 3 shows that learning by gradient descent with a  $L_2$  regularizer results in a balanced state.  
 279 Figure 4 shows that even when the network is initialized in a balanced state, without the regularizer  
 280 the network can become unbalanced if the fixed learning rate is not very small. Figure 5 shows that  
 281 the local stochastic balancing algorithm, by which neurons are randomly balanced in an asynchronous  
 282 fashion, always converges to the same (unique) global balanced state.

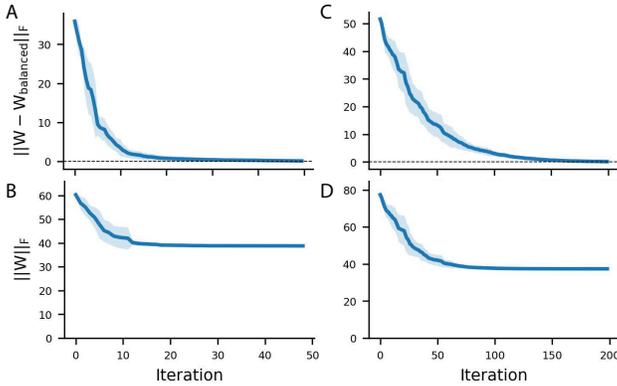


Figure 5: **Stochastic balancing converges to a unique global balanced state** (A-B) Simulations use a deep fully connected auto-encoder trained on the MNIST dataset. (C-D) Simulations use a deep locally connected network trained on the CFAR10 dataset. (A,C) The weights of the network are initialized randomly and saved. The stochastic balancing algorithm is applied and the resulting balanced weights are denoted by  $W_{balanced}$ . The stochastic balancing algorithm is applied 1,000 different times. In all repetitions, the weights converge to the same value  $W_{balanced}$ . (B,D) Stochastic balancing decreases the norm of the weights. (A-D) Shaded areas correspond to one standard deviation around the mean.

## 283 7 Conclusion

284 While the theory of neural synaptic balance is a mathematical theory that stands on its own, it is  
 285 worth considering some of its possible consequences and applications, at the theoretical, algorithmic,  
 286 biological, and neuromorphic hardware levels. At the theory level, for instance, it suggests extending  
 287 theorems obtained with ReLU neurons to BiLU neurons, using balance ideas to study learning in  
 288 *linear* regularized networks, and using the manifolds of equivalent weights to study issues of over-  
 289 parameterization (e.g. the data needs only to specify the balanced state, not the entire equivalence  
 290 class). At the algorithmic level, balancing algorithms could be used for instance to balance networks  
 291 at any stage of learning, including at the beginning, and as an alternative way to regularize networks.  
 292 Finally, because scaling and balancing are local operations, they are potentially of interest in physical,  
 293 as opposed to digitally-simulated, neural networks. In particular, it would be interesting to know if  
 294 some notion of balance applies to biological neurons. Unfortunately, current recording technologies  
 295 do not allow the measurement of all incoming and outgoing synapses of a neuron. Perhaps some  
 296 approximation could be obtained statistically and at the population level, or perhaps approximate  
 297 measurements could be carried in very simple networks (e.g. *C. elegans*) or using neurons in culture.  
 298 Finally, in neuromorphic hardware, the balance could be relevant for training spiking neural networks  
 299 with low energy consumption [Sorbaro et al., 2020, Rueckauer et al., 2017]). In particular, ReLU  
 300 scaling can influence the number of spikes generated in each layer and the average energy consumption  
 301 at each layer. Similarly, in memristor networks [Ivanov et al., 2022, Liang and Wong, 2000]),  $L_2$   
 302 minimization is directly connected to power consumption. Moreover, the issue of the limited  
 303 conductivity range of memristors is mentioned in Ivanov et al. [2022] and in Ji et al. [2016] Therefore,  
 304 a local algorithm to reduce the norm of the weights could help mitigate this issue as well.

305 The theory of neural synaptic balance explains some basic findings regarding  $L_2$  balance in feedfor-  
 306 ward networks of ReLU neurons and extends them in several directions. The first direction is the  
 307 extension to BiLU and other activation functions (BiPU). The second direction is the extension to  
 308 more general regularizers, including all  $L_p$  ( $p > 0$ ) regularizers. The third direction is the extension to  
 309 non-layered architectures, recurrent architectures, convolutional architectures, as well as architectures  
 310 with mixed activation functions. The theory is based on two local neuronal operations: scaling  
 311 which is commutative, and balancing which is not commutative. Finally, and most importantly, given  
 312 any initial set of weights, when local balancing operations are applied in a stochastic or determin-  
 313 istic manner, global order always emerges through the convergence of the balancing algorithm to  
 314 the same unique set of balanced weights. The reason for this convergence is the existence of an  
 315 underlying convex optimization problem where the relevant variables are constrained to a linear,  
 316 only architecture-dependent, manifold. Scaling and balancing operations are local and thus may  
 317 have applications in physical, non-digitally simulated, neural networks where the emergence of  
 318 global order from local operations may lead to better operating characteristics and lower energy  
 319 consumption.

## 320 **References**

- 321 P. Baldi. *Deep Learning in Science*. Cambridge University Press, Cambridge, UK, 2021.
- 322 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal*  
323 *Processing Magazine*, 29(6):141–142, 2012.
- 324 Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous  
325 models: Layers are automatically balanced. *Advances in Neural Information Processing Systems*,  
326 31, 2018.
- 327 Rachel E Field, James A D’amour, Robin Tremblay, Christoph Miehl, Bernardo Rudy, Julijana  
328 Gjorgjieva, and Robert C Froemke. Heterosynaptic plasticity determines the set point for cortical  
329 excitatory-inhibitory balance. *Neuron*, 106(5):842–854, 2020.
- 330 Robert C Froemke. Plasticity of cortical excitatory-inhibitory balance. *Annual review of neuroscience*,  
331 38:195–219, 2015.
- 332 Oliver D Howes and Ekaterina Shatalina. Integrating the neurodevelopmental and dopamine hypothe-  
333 ses of schizophrenia and the role of cortical excitation-inhibition balance. *Biological psychiatry*,  
334 2022.
- 335 Dmitry Ivanov, Aleksandr Chezhegov, Mikhail Kiselev, Andrey Grunin, and Denis Larionov. Neuro-  
336 morphic artificial intelligence systems. *Frontiers in Neuroscience*, 16:1513, 2022.
- 337 Yu Ji, YouHui Zhang, ShuangChen Li, Ping Chi, CiHang Jiang, Peng Qu, Yuan Xie, and WenGuang  
338 Chen. Neutrams: Neural network transformation and co-design under neuromorphic hardware  
339 constraints. In *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture*  
340 *(MICRO)*, pages 1–13. IEEE, 2016.
- 341 Dongshin Kim and Jang-Sik Lee. Neurotransmitter-induced excitatory and inhibitory functions in  
342 artificial synapses. *Advanced Functional Materials*, 32(21):2200497, 2022.
- 343 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- 344 Faming Liang and Wing Hung Wong. Evolutionary monte carlo: Applications to cp model sampling  
345 and change point problem. *STATISTICA SINICA*, 10:317–342, 2000.
- 346 Behnam Neyshabur, Ryota Tomioka, Ruslan Salakhutdinov, and Nathan Srebro. Data-dependent path  
347 normalization in neural networks. *arXiv preprint arXiv:1511.06747*, 2015.
- 348 Bodo Rueckauer, Iulia-Alexandra Lungu, Yuhuang Hu, Michael Pfeiffer, and Shih-Chii Liu. Conver-  
349 sion of continuous-valued deep networks to efficient event-driven networks for image classification.  
350 *Frontiers in neuroscience*, 11:294078, 2017.
- 351 Farshad Shirani and Hannah Choi. On the physiological and structural contributors to the dynamic  
352 balance of excitation and inhibition in local cortical networks. *bioRxiv*, pages 2023–01, 2023.
- 353 Martino Sorbaro, Qian Liu, Massimo Bortone, and Sadique Sheik. Optimizing the energy consump-  
354 tion of spiking neural networks for neuromorphic applications. *Frontiers in neuroscience*, 14:662,  
355 2020.
- 356 Christopher H Stock, Sarah E Harvey, Samuel A Ocko, and Surya Ganguli. Synaptic balancing: A  
357 biologically plausible local learning rule that provably increases neural network noise robustness  
358 without sacrificing task performance. *PLOS Computational Biology*, 18(9):e1010418, 2022.
- 359 A. Tavakoli, F. Agostinelli, and P. Baldi. SPLASH: Learnable activation functions for improving  
360 accuracy and adversarial robustness. *Neural Networks*, 140:1–12, 2021. Also: arXiv:2006.08947.

## 361 Appendix

### 362 A Homogeneous and BiLU Activation Functions

363 In this section, we generalize the basic example of the introduction from the standpoint of the  
364 activation functions. In particular, we consider homogeneous activation functions (defined below).  
365 The importance of homogeneity has been previously identified in somewhat different contexts  
366 [Neyshabur et al. \[2015\]](#). Intuitively, homogeneity is a form of linearity with respect to weight scaling  
367 and thus it is useful to motivate the concept of homogeneous activation functions by looking at other  
368 notions of linearity for activation functions. This will also be useful for Section E where even more  
369 general classes of activation functions are considered.

#### 370 A.1 Additive Activation Functions

371 **Definition A.1.** A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is additively linear if and only if  
372  $f(x + y) = f(x) + f(y)$  for any real numbers  $x$  and  $y$ .

373 **Proposition A.2.** The class of additively linear activation functions is exactly equal to the class of  
374 linear activation functions, i.e., activation functions of the form  $f(x) = ax$ .

375 *Proof.* Obviously linear activation functions are additively linear. Conversely, if  $f$  is additively linear,  
376 the following three properties are true:

377 (1) One must have:  $f(nx) = nf(x)$  and  $f(x/n) = f(x)/n$  for any  $x \in \mathbb{R}$  and any  $n \in \mathbb{N}$ . As a  
378 result,  $f(n/m) = nf(1)/m$  for any integers  $n$  and  $m$  ( $m \neq 0$ ).

379 (2) Furthermore,  $f(0 + 0) = f(0) + f(0)$  which implies:  $f(0) = 0$ .

380 (3) And thus  $f(x - x) = f(x) + f(-x) = 0$ , which in turn implies that  $f(-x) = -f(x)$ .

381 From these properties, it is easy to see that  $f$  must be continuous, with  $f(x) = xf(1)$ , and thus  $f$   
382 must be linear.  $\square$

#### 383 A.2 Multiplicative Activation Functions

384 **Definition A.3.** A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is multiplicative if and only if  $f(xy) =$   
385  $f(x)f(y)$  for any real numbers  $x$  and  $y$ .

386 **Proposition A.4.** The class of continuous multiplicative activation functions is exactly equal to the  
387 class of functions comprising the functions:  $f(x) = 0$  for every  $x$ ,  $f(x) = 1$  for every  $x$ , and all the  
388 even and odd functions satisfying  $f(x) = x^c$  for  $x \geq 0$ , where  $c$  is any constant in  $\mathbb{R}$ .

389 *Proof.* It is easy to check the functions described in the proposition are multiplicative. Conversely,  
390 assume  $f$  is multiplicative. For both  $x = 0$  and  $x = 1$ , we must have  $f(x) = f(xx) = f(x)f(x)$  and  
391 thus  $f(0)$  is either 0 or 1, and similarly for  $f(1)$ . If  $f(1) = 0$ , then for any  $x$  we must have  $f(x) = 0$   
392 because:  $f(x) = f(1x) = f(1)f(x) = 0$ . Likewise, if  $f(0) = 1$ , then for any  $x$  we must have  
393  $f(x) = 1$  because:  $1 = f(0) = f(0x) = f(0)f(x) = f(x)$ . Thus, in the rest of the proof, we can  
394 assume that  $f(0) = 0$  and  $f(1) = 1$ . By induction, it is easy to see that for any  $x \geq 0$  we must have:  
395  $f(x^n) = f(x)^n$  and  $f(x^{1/n}) = (f(x))^{1/n}$  for any integer (positive or negative). As a result, for any  
396  $x \in \mathbb{R}$  and any integers  $n$  and  $m$  we must have:  $f(x^{n/m}) = f(x)^{n/m}$ . By continuity this implies  
397 that for any  $x \geq 0$  and any  $r \in \mathbb{R}$ , we must have:  $f(x^r) = f(x)^r$ . Now there is some constant  $c$  such  
398 that:  $f(e) = e^c$ . And thus, for any  $x > 0$ ,  $f(x) = f(e^{\log x}) = [f(e)]^{\log x} = e^{c \log x} = x^c$ . To address  
399 negative values of  $x$ , note that we must have  $f[(-1)(-1) = f(1) = 1f(-1)^2]$ . Thus,  $f(-1)$  is either  
400 equal to 1 or to -1. Since for any  $x > 0$  we have  $f(-x) = f(-1)f(x)$ , we see that if  $f(-1) = 1$   
401 the function must be even ( $f(-x) = f(x) = x^c$ ), and if  $f(-1) = -1$  the function must be odd  
402 ( $f(-x) = -f(x)$ ).  $\square$

403 We will return to multiplicative activation function in a later section.

404 **A.3 Linearly Scalable Activation Functions**

405 **Definition A.5.** A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is linearly scalable if and only if  
406  $f(\lambda x) = \lambda f(x)$  for every  $\lambda \in \mathbb{R}$ .

407 **Proposition A.6.** The class of linearly scalable activation functions is exactly equal to the class of  
408 linear activation functions, i.e., activation functions of the form  $f(x) = ax$ .

409 *Proof.* Obviously, linear activation functions are linearly scalable. For the converse, if  $f$  is linearly  
410 multiplicative we must have  $f(\lambda x) = \lambda f(x) = x f(\lambda)$  for any  $x$  and any  $\lambda$ . By taking  $\lambda = 1$ , we get  
411  $f(x) = f(1)x$  and thus  $f$  is linear.  $\square$

412 Thus the concepts of linearly additive or linearly scalable activation function are of limited interest  
413 since both of them are equivalent to the concept of linear activation function. A more interesting  
414 class is obtained if we consider linearly scalable activation functions, where the scaling factor  $\lambda$  is  
415 constrained to be positive ( $\lambda > 0$ ), also called homogeneous functions.

416 **A.4 Homogeneous Activation Functions**

417 **Definition A.7.** (Homogeneous) A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is homogeneous if and  
418 only if:  $f(\lambda x) = \lambda f(x)$  for every  $\lambda \in \mathbb{R}$  with  $\lambda > 0$ .

419 *Remark A.8.* Note that if  $f$  is homogeneous,  $f(\lambda 0) = \lambda f(0) = f(0)$  for any  $\lambda > 0$  and thus  
420  $f(0) = 0$ . Thus it makes no difference in the definition of homogeneous if we set  $\lambda \geq 0$  instead of  
421  $\lambda > 0$ .

422 *Remark A.9.* Clearly, linear activation functions are homogeneous. However, there exists also  
423 homogeneous functions that are non-linear, such as ReLU or leaky ReLU activation functions.

424 We now provide a full characterization of the class of homogeneous activation functions.

425 **A.5 BiLU Activation Functions**

426 We first define a new class of activation functions, corresponding to bilinear units (BiLU), consisting  
427 of two half-lines meeting at the origin. This class contains all the linear functions, as well as the  
428 ReLU and leaky ReLU functions, and many other functions.

429 **Definition A.10.** (BiLU) A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is bilinear (BiLU) if and only if  
430  $f(x) = ax$  when  $x < 0$ , and  $f(x) = bx$  when  $x \geq 0$ , for some fixed parameters  $a$  and  $b$  in  $\mathbb{R}$ .

431 These include linear units ( $a = b$ ), ReLU units ( $a = 0, b = 1$ ), leaky ReLU ( $a = \epsilon; b = 1$ ) units,  
432 and symmetric linear units ( $a = -b$ ), all of which can also be viewed as special cases of piece-wise  
433 linear units [Tavakoli et al. \[2021\]](#), with a single hinge. One advantage of ReLU and more generally  
434 BiLU neurons, which is very important during backpropagation learning, is that their derivative is  
435 very simple and can only take one of two values ( $a$  or  $b$ ).

436 **Proposition A.11.** A neuronal activation function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is homogeneous if and only if it is a  
437 BiLU activation function.

438 *Proof.* Every function in BiLU is clearly homogeneous. Conversely, any homogeneous function  $f$   
439 must satisfy: (1)  $f(0x) = 0f(x) = f(0) = 0$ ; (2)  $f(x) = f(1x) = f(1)x$  for any positive  $x$ ; and (3)  
440  $f(x) = f(-u) = f(-1)u = -f(-1)x$  for any negative  $x$ . Thus  $f$  is in BiLU with  $a = -f(-1)$   
441 and  $b = f(1)$ .  $\square$

442 In Appendix A, we provide a simple proof that networks of BiLU neurons, even with a single  
443 hidden layer, have universal approximation properties. In the next two sections, we introduce two  
444 fundamental neuronal operations, scaling and balancing, that can be applied to the incoming and  
445 outgoing synaptic weights of neurons with BiLU activation functions.

## 446 B Scaling

447 **Definition B.1.** (*Scaling*) For any BiLU neuron  $i$  in network and any  $\lambda > 0$ , we let  $S_\lambda(i)$  denote the  
 448 synaptic scaling operation by which the incoming connection weights of neuron  $i$  are multiplied by  $\lambda$   
 449 and the outgoing connection weights of neuron  $i$  are divided by  $\lambda$ .

450 Note that because of the homogeneous property, the scaling operation does not change how neuron  $i$   
 451 affects the rest of the network. In particular, the input-output function of the overall network remains  
 452 unchanged after scaling neuron  $i$  by any  $\lambda > 0$ . Note also that scaling always preserves the sign of  
 453 the synaptic weights to which it is applied, and the scaling operation can never convert a non-zero  
 454 synaptic weight into a zero synaptic weight, or vice versa.

455 As usual, the bias is treated here as an additional synaptic weight emanating from a unit clamped to  
 456 the value one. Thus scaling is applied to the bias.

457 **Proposition B.2.** (*Commutativity of Scaling*) Scaling operations applied to any pair of BiLU neurons  
 458  $i$  and  $j$  in a neural network commute:  $S_\lambda(i)S_\mu(j) = S_\mu(j)S_\lambda(i)$ , in the sense that the resulting  
 459 network weights are the same, regardless of the order in which the scaling operations are applied.  
 460 Furthermore, for any BiLU neuron  $i$ :  $S_\lambda(i)S_\mu(i) = S_\mu(i)S_\lambda(i) = S_{\lambda\mu}(i)$ .

461 This is obvious. As a result, any set  $I$  of BiLU neurons in a network can be scaled simultaneously or  
 462 in any sequential order while leading to the same final configuration of synaptic weights. If we denote  
 463 by  $1, 2, \dots, n$  the neurons in  $I$ , we can for instance write:  $\prod_{i \in I} S_{\lambda_i}(i) = \prod_{\sigma(i) \in I} S_{\lambda_{\sigma(i)}}(\sigma(i))$  for  
 464 any permutation  $\sigma$  of the neurons. Likewise, we can collapse operations applied to the same neuron.  
 465 For instance, we can write:  $S_5(1)S_2(2)S_3(1)S_4(2) = S_{15}(1)S_8(2) = S_8(2)S_{15}(1)$

466 **Definition B.3.** (*Coordinated Scaling*) For any set  $I$  of BiLU neurons in a network and any  $\lambda > 0$ ,  
 467 we let  $S_\lambda(I)$  denote the synaptic scaling operation by which all the neurons in  $I$  are scaled by the  
 468 same  $\lambda$ .

## 469 C Balancing

470 **Definition C.1.** (*Balancing*) Given a BiLU neuron in a network, the balancing operation  $B(i)$  is  
 471 a particular scaling operation  $B(i) = S_{\lambda^*}(i)$ , where the scaling factor  $\lambda^*$  is chosen to optimize a  
 472 particular cost function, or regularizer, associated with the incoming and outgoing weights of neuron  
 473  $i$ .

474 For now, we can imagine that this cost function is the usual  $L_2$  (least squares) regularizer, but in  
 475 the next section, we will consider more general classes of regularizers and study the corresponding  
 476 optimization process. For the  $L_2$  regularizer, as shown in the next section, this optimization process  
 477 results in a unique value of  $\lambda^*$  such that sum of the squares of the incoming weights is equal to  
 478 the sum of the squares of the outgoing weights, hence the term ‘‘balance’’. Note that obviously  
 479  $B(B(i)) = B(i)$  and that, as a special case of scaling operation, the balancing operation does not  
 480 change how neuron  $i$  contributes to the rest of the network, and thus it leaves the overall input-output  
 481 function of the network unchanged.

482 Unlike scaling operations, balancing operations in general do not commute as balancing operations  
 483 (they still commute as scaling operations). Thus, in general,  $B(i)B(j) \neq B(j)B(i)$ . This is because  
 484 if neuron  $i$  is connected to neuron  $j$ , balancing  $i$  will change the connection between  $i$  and  $j$ , and, in  
 485 turn, this will change the value of the optimal scaling constant for neuron  $j$  and vice versa. However,  
 486 if there are no non-zero connections between neuron  $i$  and neuron  $j$  then the balancing operations  
 487 commute since each balancing operation will modify a different, non-overlapping, set of weights.

488 **Definition C.2.** (*Disjoint neurons*) Two neurons  $i$  and  $j$  in a neural network are said to be disjoint if  
 489 there are no non-zero connections between  $i$  and  $j$ .

490 Thus in this case  $B(i)B(j) = S_{\lambda^*}(i)S_{\mu^*}(j) = S_{\mu^*}(j)S_{\lambda^*}(i) = B(j)B(i)$ . This can be extended to  
 491 disjoint sets of neurons.

492 **Definition C.3.** (*Disjoint Set of Neurons*) A set  $I$  of neurons is said to be disjoint if for any pair  $i$  and  
 493  $j$  of neurons in  $I$  there are no non-zero connections between  $i$  and  $j$ .

494 For example, in a layered feedforward network, all the neurons in a layer form a disjoint set, as long  
 495 as there are no intra-layer connections or, more precisely, no non-zero intra-layer connections. All

496 the neurons in a disjoint set can be balanced in any order resulting in the same final set of synaptic  
 497 weights. Thus we have:

498 **Proposition C.4.** *If we index by  $1, 2, \dots, n$  the neurons in a disjoint set  $I$  of BiLU neurons in a*  
 499 *network, we have:  $\prod_{i \in I} B(i) = \prod_{i \in I} S_{\lambda_i^*}(i) = \prod_{\sigma(i) \in I} S_{\lambda_{\sigma(i)}^*}(\sigma(i)) = \prod_{\sigma(i) \in I} B(\sigma(i))$  for any*  
 500 *permutation  $\sigma$  of the neurons.*

501 Finally, we can define the coordinated balancing of any set  $I$  of BiLU neurons (disjoint or not  
 502 disjoint).

503 **Definition C.5.** *(Coordinated Balancing) Given any set  $I$  of BiLU neurons (disjoint or not disjoint)*  
 504 *in a network, the coordinated balancing of these neurons, written as  $B_{\lambda^*}(I)$ , corresponds to the*  
 505 *coordinated scaling all the neurons in  $I$  by the same factor  $\lambda^*$ , Where  $\lambda^*$  minimizes the cost functions*  
 506 *of all the weights, incoming and outgoing, associated with all the neurons in  $I$ .*

507 *Remark C.6.* While balancing corresponds to a full optimization of the scaling operation, it is also  
 508 possible to carry a partial optimization of the scaling operation by choosing a scaling factor that  
 509 reduces the corresponding contribution to the regularizer without minimizing it.

## 510 D General Framework and Single Neuron Balance

511 In this section, we generalize the kinds of regularizer to which the notion of neuronal synaptic balance  
 512 can be applied, beyond the usual  $L_2$  regularizer and derive the corresponding balance equations.  
 513 Thus we consider a network (feedforward or recurrent) where the hidden units are BiLU units.  
 514 The visible units can be partitioned into input units and output units. For any hidden unit  $i$ , if we  
 515 multiply all its incoming weights  $IN(i)$  by some  $\lambda > 0$  and all its outgoing weights  $OUT(i)$  by  
 516  $1/\lambda$  the overall function computed by the network remains unchanged due to the BiLU homogeneity  
 517 property. In particular, if there is an error function that depends uniquely on the input-output function  
 518 being computed, this error remains unchanged by the introduction of the multiplier  $\lambda$ . However, if  
 519 there is also a regularizer  $R$  for the weights, its value is affected by  $\lambda$  and one can ask what is the  
 520 optimal value of  $\lambda$  with respect to the regularizer, and what are the properties of the resulting weights.  
 521 This approach can be applied to any regularizer. For most practical purposes, we can assume that  
 522 the regularizer is continuous in the weights (hence in  $\lambda$ ) and lower-bounded. Without any loss of  
 523 generality, we can assume that it is lower-bounded by zero. If we want the minimum value to be  
 524 achieved by some  $\lambda > 0$ , we need to add some mild condition that prevents the minimal value from  
 525 being approached as  $\lambda \rightarrow 0$ , or as  $\lambda \rightarrow +\infty$ . For instance, it is enough if there is an interval  $[a, b]$   
 526 with  $0 < a < b$  where  $R$  achieves a minimal value  $R_{min}$  and  $R \geq R_{min}$  in the intervals  $(0, a]$  and  
 527  $[b, +\infty)$ . Additional (mild) conditions must be imposed if one wants the optimal value of  $\lambda$  to be  
 528 unique, or computable in closed form (see Theorems below). Finally, we want to be able to apply the  
 529 balancing approach

530 Thus, we consider overall regularized error functions, where the regularizer is very general, as long  
 531 as it has an additive form with respect to the individual weights:

$$532 \mathcal{E}(W) = E(W) + R(W) \quad \text{with} \quad R(W) = \sum_w g_w(w) \quad (\text{D.1})$$

532 where  $W$  denotes all the weights in the network and  $E(W)$  is typically the negative log-likelihood  
 533 (LMS error in regression tasks, or cross-entropy error in classification tasks). We assume that the  $g_w$   
 534 are continuous, and lower-bounded by 0. To ensure the existence and uniqueness of minimum during  
 535 the balancing of any neuron, We will assume that each function  $g_w$  depends only on the magnitude  
 536  $|w|$  of the corresponding weight, and that  $g_w$  is monotonically increasing from 0 to  $+\infty$  ( $g_w(0) = 0$   
 537 and  $\lim_{x \rightarrow +\infty} g_w(x) = +\infty$ ). Clearly,  $L_2, L_1$  and more generally all  $L_p$  regularizers are special  
 538 cases where, for  $p > 0$ ,  $L^p$  regularization is defined by:  $R(W) = \sum_w |w|^p$ .

539 When indicated, we may require also that the functions  $g_w$  be continuously differentiable, except  
 540 perhaps at the origin in order to be able to differentiate the regularizer with respect to the  $\lambda$ 's and  
 541 derive closed form conditions for the corresponding optima. This is satisfied by all forms of  $L_p$   
 542 regularization, for  $p > 0$ .

543 *Remark D.1.* Often one introduces scalar multiplicative hyperparameters to balance the effect of  $E$   
 544 and  $R$ , for instance in the form:  $\mathcal{E} = E + \beta R$ . These cases are included in the framework above:  
 545 multipliers like  $\beta$  can easily be absorbed into the functions  $g_w$  above.

546 **Theorem D.2.** (General Balance Equation). Consider a neural network with BiLU activation  
 547 functions in all the hidden units and overall error function of the form:

$$\mathcal{E} = E(W) + R(W) \quad \text{with} \quad R(W) = \sum_w g_w(w) \quad (\text{D.2})$$

548 where each function  $g_w(w)$  is continuous, depends on the magnitude  $|w|$  alone, and grows monotonically  
 549 from  $g_w(0) = 0$  to  $g_w(+\infty) = +\infty$ . For any setting of the weights  $W$  and any hidden unit  $i$  in  
 550 the network and any  $\lambda > 0$  we can multiply the incoming weights of  $i$  by  $\lambda$  and the outgoing weights  
 551 of  $i$  by  $1/\lambda$  without changing the overall error  $E$ . Furthermore, there exists a unique value  $\lambda^*$  where  
 552 the corresponding weights  $v$  ( $v = \lambda^* w$  for incoming weights,  $v = w/\lambda^*$  for the outgoing weights)  
 553 achieve the balance equation:

$$\sum_{v \in IN(i)} g_w(v) = \sum_{w \in OUT(i)} g_w(v) \quad (\text{D.3})$$

554 *Proof.* Under the assumptions of the theorem,  $E$  is unchanged under the rescaling of the incoming and  
 555 outgoing weights of unit  $i$  due to the homogeneity property of BiLUs. Without any loss of generality,  
 556 let us assume that at the beginning:  $\sum_{w \in IN(i)} g_w(w) < \sum_{w \in OUT(i)} g_w(w)$ . As we increase  $\lambda$  from  
 557 1 to  $+\infty$ , by the assumptions on the functions  $g_w$ , the term  $\sum_{w \in IN(i)} g_w(\lambda w)$  increases continuously  
 558 from its initial value to  $+\infty$ , whereas the term  $\sum_{w \in OUT(i)} g_w(w/\lambda)$  decreases continuously from  
 559 its initial value to 0. Thus, there is a unique value  $\lambda^*$  where the balance is realized. If at the beginning  
 560  $\sum_{w \in IN(i)} g_w(w) > \sum_{w \in OUT(i)} g_w(w)$ , then the same argument is applied by decreasing  $\lambda$  from 1  
 561 to 0.  $\square$

562 *Remark D.3.* For simplicity, here and in other sections, we state the results in terms of a network of  
 563 BiLU units. However, the same principles can be applied to networks where only a subset of neurons  
 564 are in the BiLU class, simply by applying scaling and balancing operations to only those neurons.  
 565 Furthermore, not all BiLU neurons need to have the same BiLU activation function. For instance, the  
 566 results still hold for a mixed network containing both ReLU and linear units.

567 *Remark D.4.* In the setting of Theorem D.2, the balance equations do not necessarily minimize the  
 568 corresponding regularization term. This is addressed in the next theorem.

569 *Remark D.5.* Finally, zero weights ( $w = 0$ ) can be ignored entirely as they play no role in scaling or  
 570 balancing. Furthermore, if all the incoming or outgoing weights of a hidden unit were to be zero, it  
 571 could be removed entirely from the network

572 **Theorem D.6.** (Balance and Regularizer Minimization) We now consider the same setting as in  
 573 Theorem D.2, but in addition, we assume that the functions  $g_w$  are continuously differentiable, except  
 574 perhaps at the origin. Then, for any neuron, there exists at least one optimal value  $\lambda^*$  that minimizes  
 575  $R(W)$ . This value must be a solution of the consistency equation:

$$\lambda^2 \sum_{w \in IN(i)} w g'_w(\lambda w) = \sum_{w \in OUT(i)} w g'_w(w/\lambda) \quad (\text{D.4})$$

576 Once the weights are rebalanced accordingly, the new weights must satisfy the generalized balance  
 577 equation:

$$\sum_{w \in IN(i)} w g'(w) = \sum_{w \in OUT(i)} w g'(w) \quad (\text{D.5})$$

578 In particular, if  $g_w(w) = |w|^p$  for all the incoming and outgoing weights of neuron  $i$ , then the optimal  
 579 value  $\lambda^*$  is unique and equal to:

$$\lambda^* = \left( \frac{\sum_{w \in OUT(i)} |w|^p}{\sum_{w \in IN(i)} |w|^p} \right)^{1/2p} = \left( \frac{\|OUT(i)\|_p}{\|IN(i)\|_p} \right)^{1/2} \quad (\text{D.6})$$

580 The decrease  $\Delta R \geq 0$  in the value of the  $L_p$  regularizer  $R = \sum_w |w|^p$  is given by:

$$\Delta R = \left( \left( \sum_{w \in IN(i)} |w|^p \right)^{1/2} - \left( \sum_{w \in OUT(i)} |w|^p \right)^{1/2} \right)^2 \quad (\text{D.7})$$

581 After balancing neuron  $i$ , its new weights satisfy the generalized  $L_p$  balance equation:

$$\sum_{w \in IN(i)} |w|^p = \sum_{w \in OUT(i)} |w|^p \quad (\text{D.8})$$

582 *Proof.* Due to the additivity of the regularizer, the only component of the regularizer that depends on  
583  $\lambda$  has the form:

$$R(\lambda) = \sum_{w \in IN(i)} g_w(\lambda w) + \sum_{w \in OUT(i)} g_w(w/\lambda) \quad (\text{D.9})$$

584 Because of the properties of the functions  $g_w$ ,  $R_\lambda$  is continuously differentiable and strictly bounded  
585 below by 0. So it must have a minimum, as a function of  $\lambda$  where its derivative is zero. Its derivative  
586 with respect to  $\lambda$  has the form:

$$R'(\lambda) = \sum_{w \in IN(i)} w g'_w(\lambda w) + \sum_{w \in OUT(i)} (-w/\lambda^2) g'_w(w/\lambda) \quad (\text{D.10})$$

587 Setting the derivative to zero, gives:

$$\lambda^2 \sum_{w \in IN(i)} w g'_w(\lambda w) = \sum_{w \in OUT(i)} w g'_w(w/\lambda) \quad (\text{D.11})$$

588 Assuming that the left-hand side is non-zero, which is generally the case, the optimal value for  $\lambda$   
589 must satisfy:

$$\lambda = \left( \frac{\sum_{w \in OUT(i)} w g'_w(w/\lambda)}{\sum_{w \in IN(i)} w g'_w(\lambda w)} \right)^{1/2} \quad (\text{D.12})$$

590 If the regularizing function is the same for all the incoming and outgoing weights ( $g_w = g$ ), then the  
591 optimal value  $\lambda$  must satisfy:

$$\lambda = \left( \frac{\sum_{w \in OUT(i)} w g'(w/\lambda)}{\sum_{w \in IN(i)} w g'(\lambda w)} \right)^{1/2} \quad (\text{D.13})$$

592 In particular, if  $g(w) = |w|^p$  then  $g(w)$  is differentiable except possibly at 0 and  $g'(w) =$   
593  $s(w)p|w|^{p-1}$ , where  $s(w)$  denotes the sign of the weight  $w$ . Substituting in Equation D.13, the  
594 optimal rescaling  $\lambda$  must satisfy:

$$\lambda^* = \left( \frac{\sum_{w \in OUT(i)} w s(w) |w|^{p-1}}{\sum_{w \in IN(i)} w |w s(w)|^{p-1}} \right)^{1/2p} = \left( \frac{\sum_{w \in OUT(i)} |w|^p}{\sum_{w \in IN(i)} |w|^p} \right)^{1/2p} = \left( \frac{\|OUT(i)\|_p}{\|IN(i)\|_p} \right)^{1/2} \quad (\text{D.14})$$

595 At the optimum, no further balancing is possible, and thus  $\lambda^* = 1$ . Equation D.11 yields immediately  
596 the generalized balance equation to be satisfied at the optimum:

$$\sum_{w \in IN(i)} w g'(w) = \sum_{w \in OUT(i)} w g'(w) \quad (\text{D.15})$$

597 In the case of  $L_P$  regularization, it is easy to check by applying Equation D.15, or by direct calculation  
598 that:

$$\sum_{w \in IN(i)} |\lambda^* w|^p = \sum_{w \in OUT(i)} |w/\lambda^*|^p \quad (\text{D.16})$$

599 which is the generalized balance equation. Thus after balancing neuron, the weights of neuron  $i$   
600 satisfy the  $L_p$  balance (Equation D.8). The change in the value of the regularizer is given by:

$$\Delta R = \sum_{w \in IN(i)} |w|^p + \sum_{w \in OUT(i)} |w|^p - \sum_{w \in IN(i)} |\lambda^* w|^p - \sum_{w \in OUT(i)} |w/\lambda^*|^p \quad (\text{D.17})$$

601 By substituting  $\lambda^*$  by its explicit value given by Equation D.14 and collecting terms gives Equation  
602 D.7.  $\square$

603 *Remark D.7.* The monotonicity of the functions  $g_w$  is not needed to prove the first part of Theorem  
604 D.6. It is only needed to prove the uniqueness of  $\lambda^*$  in the  $L_p$  cases.

605 *Remark D.8.* Note that the same approach applies to the case where there are multiple additive  
606 regularizers. For instance with both  $L^2$  and  $L^1$  regularization, in this case the function  $f$  has the form:  
607  $g_w(w) = \alpha w^2 + \beta |w|$ . Generalized balance still applies. It also applies to the case where different  
608 regularizers are applied in different disconnected portions of the network.

609 *Remark D.9.* The balancing of a single BiLU neuron has little to do with the number of connections.  
610 It applies equally to fully connected neurons, or to sparsely connected neurons.

## 611 E Scaling and Balancing Beyond BiLU Activation Functions

612 So far we have generalized ReLU activation functions to BiLU activation functions in the context of  
613 scaling and balancing operations with positive scaling factors. While in the following sections we  
614 will continue to work with BiLU activation functions, in this section we show that the scaling and  
615 balancing operations can be extended even further to other activation functions. The section can be  
616 skipped if one prefers to progress towards the main results on stochastic balancing.

617 Given a neuron with activation function  $f(x)$ , during scaling instead of multiplying and dividing by  
618  $\lambda > 0$ , we could multiply the incoming weights by a function  $g(\lambda)$  and divide the outgoing weights  
619 by a function  $h(\lambda)$ , as long as the activation function  $f$  satisfies:

$$f(g(\lambda)x) = h(\lambda)f(x) \quad (\text{E.1})$$

620 for every  $x \in \mathbb{R}$  to ensure that the contribution of the neuron to the rest of the network remains  
621 unchanged. Note that if the activation function  $f$  satisfies Equation E.1, so does the activation  
622 function  $-f$ . In Equation E.1,  $\lambda$  does not have to be positive—we will simply assume that  $\lambda$  belongs  
623 to some open (potentially infinite) interval  $(a, b)$ . Furthermore, the functions  $g$  and  $h$  cannot be zero  
624 for  $\lambda \in (a, b)$  since they are used for scaling. It is reasonable to assume that the functions  $g$  and  $h$   
625 are continuous, and thus they must have a constant sign as  $\lambda$  varies over  $(a, b)$ .

626 Now, taking  $x = 0$  gives  $f(0) = h(\lambda)f(0)$  for every  $\lambda \in (a, b)$ , and thus either  $f(0) = 0$  or  $h(\lambda) = 1$   
627 for every  $\lambda \in (a, b)$ . The latter is not interesting and thus we can assume that the activation function  
628  $f$  satisfies  $f(0) = 0$ . Taking  $x = 1$  gives  $f(g(\lambda)) = h(\lambda)f(1)$  for every  $\lambda$  in  $(a, b)$ . For simplicity,  
629 let us assume that  $f(x) = 1$ . Then, we have:  $f(g(\lambda)) = h(\lambda)$  for every  $\lambda$ . Substituting in Equation  
630 E.1 yields:

$$f(g(\lambda)x) = f(g(\lambda))f(x) \quad (\text{E.2})$$

631 for every  $x \in \mathbb{R}$  and every  $\lambda \in (a, b)$ . This relation is essentially the same as the relation that defines  
632 multiplicative activation functions over the corresponding domain (see Proposition A.4), and thus  
633 we can identify a key family of solutions using power functions. Note that we can define a new  
634 parameter  $\mu = g(\lambda)$ , where  $\mu$  ranges also over some positive or negative interval  $I$  over which:  
635  $f(\mu x) = f(\mu)f(x)$ .

### 636 E.1 Bi-Power Units (BiPU)

637 Let us assume that  $\lambda > 0$ ,  $g(\lambda) = \lambda$  and  $h(\lambda) = \lambda^c$  for some  $c \in \mathbb{R}$ . Then the activation function  
638 must satisfy the equation:

$$f(\lambda x) = \lambda^c f(x) \quad (\text{E.3})$$

639 for any  $x \in \mathbb{R}$  and any  $\lambda > 0$ . Note that if  $f(x) = x^c$  we get a multiplicative activation function.  
 640 More generally, these functions are characterized by the following proposition.

641 **Proposition E.1.** *The set of activation functions  $f$  satisfying  $f(\lambda x) = \lambda^c f(x)$  for any  $x \in \mathbb{R}$  and  
 642 any  $\lambda > 0$  consist of the functions of the form:*

$$f(x) = \begin{cases} Cx^c & \text{if } x \geq 0 \\ Dx^c & \text{if } x < 0. \end{cases} \quad (\text{E.4})$$

643 where  $c \in \mathbb{R}$ ,  $C = f(1) \in \mathbb{R}$ , and  $D = f(-1) \in \mathbb{R}$ . We call these bi-power units (BiPU). If, in  
 644 addition, we want  $f$  to be continuous at 0, we must have either  $c > 0$ , or  $c = 0$  with  $C = D$ .

645 Given the general shape, these activations functions can be called BiPU (Bi-Power-Units). Note that  
 646 in the general case where  $c > 0$ ,  $C$  and  $D$  do not need to be equal. In particular, one of them can  
 647 be equal to zero, and the other one can be different from zero giving rise to “rectified power units”  
 648 (Figure 6).

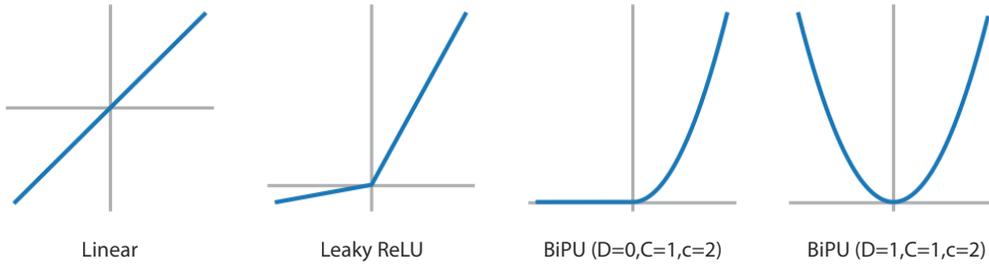


Figure 6

649 *Proof.* By taking  $x = 1$ , we get  $f(\lambda) = f(1)\lambda^c$  for any  $\lambda > 0$ . Let  $f(1) = C$ . Then we see  
 650 that for any  $x > 0$  we must have:  $f(x) = Cx^c$ . In addition, for every  $\lambda > 0$  we must have:  
 651  $f(\lambda 0) = f(0) = \lambda^c f(0)$ . So if  $c = 0$ , then  $f(x) = C = f(1)$  for  $x \geq 0$ . If  $c \neq 0$ , then  $f(0) = 0$ . In  
 652 this case, if we want the activation function to be continuous, then we see that we must have  $c \geq 0$ . So  
 653 in summary for  $x > 0$  we must have  $f(x) = f(1)x^c = Cx^c$ . For the function to be right continuous  
 654 at 0, we must have either  $f(0) = f(1) = C$  with  $c = 0$  or  $f(0) = 0$  with  $c > 0$ . We can now look  
 655 at negative values of  $x$ . By the same reasoning, we have  $f(\lambda(-1)) = f(-\lambda) = \lambda^c f(-1)$  for any  
 656  $\lambda > 0$ . Thus for any  $x < 0$  we must have:  $f(x) = f(-1)|x|^c = D|x|^c$  where  $D = f(-1)$ . Thus, if  
 657  $f$  is continuous, there are two possibilities. If  $c = 0$ , then we must have  $C = f(1) = D(f(-1))$ — and  
 658 thus  $f(x) = C$  everywhere. If  $c \neq 0$ , then continuity requires that  $c > 0$ . In this case  $f(x) = Cx^c$   
 659 for  $x \geq 0$  with  $C = f(1)$ , and  $f(x) = Dx^c$  for  $x < 0$  with  $f(-1) = D$ . In all cases, it is easy to  
 660 check directly that the resulting functions satisfy the functional equation given by Equation E.3.  $\square$

## 661 E.2 Scaling BiPU Neurons

662 A BiPU neuron can be scaled by multiplying its incoming weight by  $\lambda > 0$  and dividing its outgoing  
 663 weights by  $1/\lambda^c$ . This will not change the role of the corresponding unit in the network, and thus it  
 664 will not change the input-output function of the network.

## 665 E.3 Balancing BiPU Neurons

666 As in the case of BiLU neurons, we balance a multiplicative neuron by asking what is the optimal  
 667 scaling factor  $\lambda$  that optimizes a particular regularizer. For simplicity, here we assume that the  
 668 regularizer is in the  $L_p$  class. Then we are interested in the value of  $\lambda > 0$  that minimizes the  
 669 function:

$$\lambda^p \sum_{w \in IN} |w|^p + \frac{1}{\lambda^{pc}} \sum_{w \in OUT} |w|^p \quad (\text{E.5})$$

670 A simple calculation shows that the optimal value of  $\lambda$  is given by:

$$\lambda^* = \left( \frac{c \sum_{OUT} |w|^p}{\sum_{IN} |w|^p} \right)^{1/p(c+1)} \quad (\text{E.6})$$

671 Thus after balancing the weights, the neuron must satisfy the balance equation:

$$c \sum_{OUT} |w|^p = \sum_{IN} |w|^p \quad (\text{E.7})$$

672 in the new weights  $w$ .

673 So far, we have focused on balancing individual neurons. In the next two sections, we look at  
 674 balancing across all the units of a network. We first look at what happens to network balance when a  
 675 network is trained by gradient descent and then at what happens to network balance when individual  
 676 neurons are balanced iteratively in a regular or stochastic manner.

## 677 F Network Balance: Gradient Descent

678 A natural question is whether gradient descent (or stochastic gradient descent) applied to a network of  
 679 BiLU neurons, with or without a regularizer, converges to a balanced state of the network, where all  
 680 the BiLU neurons are balanced. So we first consider the case where there is no regularizer ( $\mathcal{E} = E$ ).  
 681 The results in Du et al. [2018] may suggest that gradient descent may converge to a balanced state. In  
 682 particular, they write that for any neuron  $i$ :

$$\frac{d}{dt} \left( \sum_{w \in IN(i)} w^2 - \sum_{w \in OUT(i)} w^2 \right) = 0 \quad (\text{F.1})$$

683 Thus the gradient flow exactly preserves the difference between the  $L_2$  cost of the incoming and  
 684 outgoing weights or, in other words, the derivative of the  $L_2$  balance *deficit* is zero. Thus if one were  
 685 to start from a balanced state and use an infinitesimally small learning rate one ought to stay in a  
 686 balanced state at all times.

687 However, it must be noted that this result was derived for the  $L_2$  metric only, and thus would not  
 688 cover other  $L_p$  forms of balance. Furthermore, it requires an infinitesimally small learning rate. In  
 689 practice, when any standard learning rate is applied, we find that gradient descent does *not* converge  
 690 to a balanced state (Figure 1). However, things are different when a regularizer term is included in  
 691 the error functions as described in the following theorem.

692 **Theorem F.1.** *Gradient descent in a network of BiLU units with error function  $\mathcal{E} = E + R$  where  $R$   
 693 has the properties described in Theorem D.6 (including all  $L_p$ ) must converge to a balanced state,  
 694 where every BiLU neuron is balanced.*

695 *Proof.* By contradiction, suppose that gradient descent converges to a state that is unbalanced and  
 696 where the gradient with respect to all the weights is zero. Then there is at least one unbalanced neuron  
 697 in the network. We can then multiply the incoming weights of such a neuron by  $\lambda$  and the outgoing  
 698 weights by  $1/\lambda$  as in the previous section without changing the value of  $E$ . Since the neuron is not in  
 699 balance, we can move  $\lambda$  infinitesimally so as to reduce  $R$ , and hence  $\mathcal{E}$ . But this contradicts the fact  
 700 that the gradient is zero.  $\square$

701 *Remark F.2.* In practice, in the case of stochastic gradient descent applied to  $E + R$ , at the end of  
 702 learning the algorithm may hover around a balanced state. If the state reached by the stochastic  
 703 gradient descent procedure is not approximately balanced, then learning ought to continue. In other  
 704 words, the degree of balance could be used to monitor whether learning has converged or not. Balance  
 705 is a necessary, but not sufficient, condition for being at the optimum.

706 *Remark F.3.* If early stopping is being used to control overfitting, there is no reason for the stopping  
 707 state to be balanced. However, the balancing algorithms described in the next section could be used  
 708 to balance this state.

## 709 G Network Balance: Stochastic or Deterministic Balancing Algorithms

710 In this section, we look at balancing algorithms where, starting from an initial weight configuration  
711  $W$ , the BiLU neurons of a network are balanced iteratively according to some deterministic or  
712 stochastic schedule that periodically visits all the neurons. We can also include algorithms where  
713 neurons are partitioned into groups (e.g. neuronal layers) and neurons in each group are balanced  
714 together.

### 715 G.1 Basic Stochastic Balancing

716 The most interesting algorithm is when the BiLU neurons of a network are iteratively balanced  
717 in a purely stochastic manner. This algorithm is particularly attractive from the standpoint of  
718 physically implemented neural networks because the balancing algorithm is local and the updates  
719 occur randomly without the need for any kind of central coordination. As we shall see in the following  
720 section, the random local operations remarkably lead to a unique form of global order. The proof  
721 for the stochastic case extends immediately to the deterministic case, where the BiLU neurons are  
722 updated in a deterministic fashion, for instance by repeatedly cycling through them according to  
723 some fixed order.

### 724 G.2 Subset Balancing (Independent or Tied)

725 It is also possible to partition the BiLU neurons into non-overlapping subsets of neurons, and then  
726 balance each subset, especially when the neurons in each subset are disjoint of each other. In this  
727 case, one can balance all the neurons in a given subset, and repeat this subset-balancing operation  
728 subset-by-subset, again in a deterministic or stochastic manner. Because the BiLU neurons in each  
729 subset are disjoint, it does not matter whether the neurons in a given subset are updated synchronously  
730 or sequentially (and in which order). Since the neurons are balanced independently of each other,  
731 this can be called independent subset balancing. For example, in a layered feedforward network with  
732 no lateral connections, each layer corresponds to a subset of disjoint neurons. The incoming and  
733 outgoing connections of each neuron are distinct from the incoming and outgoing connections of  
734 any other neuron in the layer, and thus the balancing operation of any neuron in the layer does not  
735 interfere with the balancing operation of any other neuron in the same layer. So this corresponds to  
736 independent layer balancing,

737 As a side note, balancing a layer  $h$ , may disrupt the balance of layer  $h + 1$ . However, balancing  
738 layers  $h$  and  $h + 2$  (or any other layer further apart) can be done without interference of the balancing  
739 processes. This suggests also an alternating balancing scheme, where one alternatively balances all  
740 the odd-numbered layers, and all the evenly-numbered layers.

741 Yet another variation is when the neurons in a disjoint subset are tied to each other in the sense that  
742 they must all share the same scaling factor  $\lambda$ . In this case, balancing the subset requires finding the  
743 optimal  $\lambda$  for the entire subset, as opposed to finding the optimal  $\lambda$  for each neuron in the subset.  
744 Since the neurons are balanced in a coordinated or tied fashion, this can be called coordinated or tied  
745 subset balancing. For example, tied layer balancing must use the same  $\lambda$  for all the neurons in a given  
746 layer. It is easy to see that this approach leads to layer synaptic balance which has the form (for an  
747  $L_p$  regularizer):

$$\sum_i \sum_{w \in IN(i)} |w|^p = \sum_i \sum_{w \in OUT(i)} |w|^p \quad (\text{G.1})$$

748 where  $i$  runs over all the neurons in the layer. This does *not* necessarily imply that each neuron  
749 in the layer is individually balanced. Thus neuronal balance for every neuron in a layer implies  
750 layer balance, but the converse is not true. Independent layer balancing will lead to layer balance.  
751 Coordinated layer balancing will lead to layer balance, but not necessarily to neuronal balance of  
752 each neuron in the layer. Layer-wise balancing, independent or tied, can be applied to all the layers  
753 and in a deterministic (e.g. sequential) or stochastic manner. Again the proof given in the next section  
754 for the basic stochastic algorithm can easily be applied to these cases (see also Appendix B).

755 **G.3 Remarks about Weight Sharing and Convolutional Neural Networks**

756 Suppose that two connections share the same weight so that we must have:  $w_{ij} = w_{kl}$  at all times.  
 757 In general, when the balancing algorithm is applied to neuron  $i$  or  $j$ , the weight  $w_{ij}$  will change  
 758 and the same change must be applied to  $w_{kl}$ . The latter may disrupt the balance of neuron  $k$  or  $l$ .  
 759 Furthermore, this may not lead to a decrease in the overall value of the regularizer  $R$ .

760 The case of convolutional networks is somewhat special, since *all* the incoming weights of the  
 761 neurons sharing the same convolutional kernel are shared. However, in general, the outgoing weights  
 762 are not shared. Furthermore, certain operations like max-pooling are not homogeneous. So if one  
 763 trains a CNN with  $E$  alone, or even with  $E + R$ , one should not expect any kind of balance to emerge  
 764 in the convolution units. However, all the other BiLU units in the network should become balanced  
 765 by the same argument used for gradient descent above. The balancing algorithm applied to individual  
 766 neurons, or the independent layer balancing algorithm, will not balance individual neurons sharing  
 767 the same convolution kernel. The only balancing algorithm that could lead to some convolution layer  
 768 balance, but not to individual neuronal balance, is the coordinated layer balancing, where the same  $\lambda$   
 769 is used for all the neurons in the same convolution layer, provided that their activation functions are  
 770 BiLU functions.

771 We can now study the convergence properties of balancing algorithms.

772 **H Convergence of Balancing Algorithms**

773 We now consider the basic stochastic balancing algorithm, where BiLU neurons are iteratively and  
 774 stochastically balanced. It is essential to note that balancing a neuron  $j$  may break the balance of  
 775 another neuron  $i$  to which  $j$  is connected. Thus convergence of iterated balancing is not obvious.  
 776 There are three key questions to be addressed for the basic stochastic algorithm, as well as all the  
 777 other balancing variations. First, does the value of the regularizer converge to a finite value? Second,  
 778 do the weights themselves converge to fixed finite values representing a balanced state for the entire  
 779 network? And third, if the weights converge, do they always converge to the same values, irrespective  
 780 of the order in which the units are being balanced? In other words, given an initial state  $W$  for the  
 781 network, is there a unique corresponding balanced state, with the same input-output functionalities?

782 **H.1 Notation and Key Questions**

783 For simplicity, we use a continuous time notation. After a certain time  $t$  each neuron has been  
 784 balanced a certain number of times. While the balancing operations are not commutative as balancing  
 785 operations, they are commutative as scaling operations. Thus we can reorder the scaling operations  
 786 and group them neuron by neuron so that, for instance, neuron  $i$  has been scaled by the sequence of  
 787 scaling operations:

$$S_{\lambda_1^*}(i)S_{\lambda_2^*}(i) \dots S_{\lambda_{n_{it}}^*}(i) = S_{\Lambda_i(t)}(i) \quad (\text{H.1})$$

788 where  $n_{it}$  corresponds to the count of the last update of neuron  $i$  prior to time  $t$ , and:

$$\Lambda_i(t) = \prod_{1 \leq n \leq n_{it}} \lambda_n^*(i) \quad (\text{H.2})$$

789 For the input and output units, we can consider that their balancing coefficients  $\lambda^*$  are always equal  
 790 to 1 (at all times) and therefore  $\Lambda_i(t) = 1$  for any visible unit  $i$ .

791 Thus, we first want to know if  $R$  converges. Second, we want to know if the weights converge. This  
 792 question can be split into two sub-questions: (1) Do the balancing factors  $\lambda_n^*(i)$  converge to a limit as  
 793 time goes to infinity? Even if the  $\lambda_n^*(i)$ 's converge to a limit, this does not imply that the weights of  
 794 the network converge to a limit. After a time  $t$ , the weight  $w_{ij}(t)$  between neuron  $j$  and neuron  $i$  has  
 795 the value  $w_{ij}\Lambda_i(t)/\Lambda_j(t)$ , where  $w_{ij} = w_{ij}(0)$  is the value of the weight at the start of the stochastic  
 796 balancing algorithm. Thus: (2) Do the quantities  $\Lambda_i(t)$  converge to finite values, different from 0?  
 797 And third, if the weights converge to finite values different from 0, are these values unique or not, i.e.  
 798 do they depend on the details of the stochastic updates or not? These questions are answered by the  
 799 following main theorem..

## 800 H.2 Convergence of the Basic Stochastic Balancing Algorithm to a Unique Optimum

801 **Theorem H.1.** (Convergence of Stochastic Balancing) Consider a network of BiLU neurons with an  
 802 error function  $\mathcal{E}(W) = E(W) + R(W)$  where  $R$  satisfies the conditions of Theorem D.2 including all  
 803  $L_p$  ( $p > 0$ ). Let  $W$  denote the initial weights. When the neuronal stochastic balancing algorithm is  
 804 applied throughout the network so that every neuron is visited from time to time, then  $E(W)$  remains  
 805 unchanged but  $R(W)$  must converge to some finite value that is less or equal to the initial value,  
 806 strictly less if the initial weights are not balanced. In addition, for every neuron  $i$ ,  $\lambda_i^*(t) \rightarrow 1$  and  
 807  $\Lambda_i(t) \rightarrow \Lambda_i$  as  $t \rightarrow \infty$ , where  $\Lambda_i$  is finite and  $\Lambda_i > 0$  for every  $i$ . As a result, the weights themselves  
 808 must converge to a limit  $W'$  which is globally balanced, with  $E(W) = E(W')$  and  $R(W) \geq R(W')$ ,  
 809 and with equality if only if  $W$  is already balanced. Finally,  $W'$  is unique as it corresponds to the  
 810 solution of a strictly convex optimization problem in the variables  $L_{ij} = \log(\Lambda_i/\Lambda_j)$  with linear  
 811 constraints of the form  $\sum_{\pi} L_{ij} = 0$  along any path  $\pi$  joining an input unit to an output unit and along  
 812 any directed cycle (for recurrent networks). Stochastic balancing projects to stochastic trajectories in  
 813 the linear manifold that run from the origin to the unique optimal configuration.

814 *Proof.* Each individual balancing operation leaves  $E(W)$  unchanged because the BiLU neurons are  
 815 homogeneous. Furthermore, each balancing operation reduces the regularization error  $R(W)$ , or  
 816 leaves it unchanged. Since the regularizer is lower-bounded by zero, the value of the regularizer must  
 817 approach a limit as the stochastic updates are being applied.

818 For the second question, when neuron  $i$  is balanced at some step, we know that the regularizer  $R$   
 819 decreases by:

$$\Delta R = \left( \left( \sum_{w \in IN(i)} |w|^p \right)^{1/2} - \left( \sum_{w \in OUT(i)} |w|^p \right)^{1/2} \right)^2 \quad (\text{H.3})$$

820 If the convergence were to occur in a finite number of steps, then the coefficients  $\lambda_i^*(t)$  must become  
 821 equal and constant to 1 and the result is obvious. So we can focus on the case where the convergence  
 822 does not occur in a finite number of steps (indeed this is the main scenario, as we shall see at the end  
 823 of the proof). Since  $\Delta R \rightarrow 0$ , we must have:

$$\sum_{w \in IN(i)} |w|^p \rightarrow \sum_{w \in OUT(i)} |w|^p \quad (\text{H.4})$$

824 But from the expression for  $\lambda^*$  (Equation D.14), this implies that for every  $i$ ,  $\lambda_n^*(i) \rightarrow 1$  as time  
 825 increases ( $n \rightarrow \infty$ ). This alone is not sufficient to prove that  $\Lambda_i(t)$  converges for every  $i$  as  $t \rightarrow \infty$ .  
 826 However, it is easy to see that  $\Lambda_i(t)$  cannot contain a sub-sequence that approaches 0 or  $\infty$  (Figure 7).  
 827 Furthermore, not only  $\Delta R$  converges to 0, but the series  $\sum \Delta R$  is convergent. This shows that, for  
 828 every  $i$ ,  $\Delta_i(t)$  must converge to a finite, non-zero value  $\Delta_i$ . Therefore all the weights must converge  
 829 to fixed values given by  $w_{ij}(0)\Lambda_i/\Lambda_j$ .

830 Finally, we prove that given an initial set of weights  $W$ , the final balanced state is unique and  
 831 independent of the order of the balancing operations. The coefficients  $\Lambda_i$  corresponding to a globally  
 832 balanced state must be solutions of the following optimization problem:

$$\min_{\Lambda} R(\Lambda) = \sum_{ij} \left| \frac{\Lambda_i}{\Lambda_j} w_{ij} \right|^p \quad (\text{H.5})$$

833 under the simple constraints:  $\Lambda_i > 0$  for all the BiLU hidden units, and  $\Lambda_i = 1$  for all the visible (input  
 834 and output) units. In this form, the problem is not convex. Introducing new variables  $M_j = 1/\Lambda_j$   
 835 is not sufficient to render the problem convex. Using variables  $M_{ij} = \Lambda_i/\Lambda_j$  is better, but still  
 836 problematic for  $0 < p \leq 1$ . However, let us instead introduce the new variables  $L_{ij} = \log(\Lambda_i/\Lambda_j)$ .  
 837 These are well defined since we know that  $\Lambda_i/\Lambda_j > 0$ . The objective now becomes:

$$\min R(L) = \sum_{ij} |e^{L_{ij}} w_{ij}|^p = \sum_{ij} e^{pL_{ij}} |w_{ij}|^p \quad (\text{H.6})$$

838 This objective is strictly convex in the variables  $L_{ij}$ , as a sum of strictly convex functions (exponen-  
 839 tials). However, to show that it is a convex optimization problem we need to study the constraints on

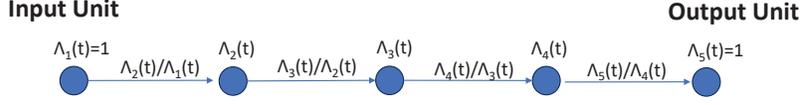


Figure 7: A path with three hidden BiLU units connecting one input unit to one output unit. During the application of the stochastic balancing algorithm, at time  $t$  each unit  $i$  has a cumulative scaling factor  $\Lambda_i(t)$ , and each directed edge from unit  $j$  to unit  $i$  has a scaling factor  $M_{ij}(t) = \Lambda_i(t)/\Lambda_j(t)$ . The  $\lambda_i(t)$  must remain within a finite closed interval away from 0 and infinity. To see this, imagine for instance that there is a subsequence of  $\Lambda_3(t)$  that approaches 0. Then there must be a corresponding subsequence of  $\Lambda_4(t)$  that approaches 0, or else the contribution of the weight  $w_{43}\Lambda_4(t)/\Lambda_3(t)$  to the regularizer would go to infinity. But then, as we reach the output layer, the contribution of the last weight  $w_{54}\Lambda_5(t)/\Lambda_4(t)$  to the regularizer goes to infinity because  $\Lambda_5(t)$  is fixed to 1 and cannot compensate for the small values of  $\Lambda_4(t)$ . And similarly, if there is a subsequence of  $\Lambda_3(t)$  going to infinity, we obtain a contradiction by propagating its effect towards the input layer.

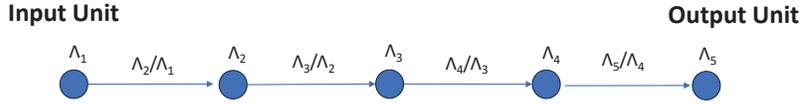


Figure 8: A path with five units. After the stochastic balancing algorithm has converged, each unit  $i$  has a scaling factor  $\Lambda_i$ , and each directed edge from unit  $j$  to unit  $i$  has a scaling factor  $M_{ij} = \Lambda_i/\Lambda_j$ . The products of the  $M_{ij}$ 's along the path is given by:  $\frac{\Lambda_2}{\Lambda_1} \frac{\Lambda_3}{\Lambda_2} \frac{\Lambda_4}{\Lambda_3} \frac{\Lambda_5}{\Lambda_4} = \frac{\Lambda_5}{\Lambda_1}$ . Accordingly, if we sum the variables  $L_{ij} = \log M_{ij}$  along the directed path, we get  $L_{21} + L_{32} + L_{43} + L_{54} = \log \Lambda_5 - \log \Lambda_1$ . In particular, if unit 1 is an input unit and unit 5 is an output unit, we must have  $\Lambda_1 = \Lambda_5 = 1$  and thus:  $L_{21} + L_{32} + L_{43} + L_{54} = 0$ . Likewise, in the case of a directed cycle where unit 1 and unit 5 are the same, we must have:  $L_{21} + L_{32} + L_{43} + L_{54} + L_{15} = 0$ .

840 the variables  $L_{ij}$ . From the set of  $\Lambda_i$ 's it is easy to construct a unique set of  $L_{ij}$ . However what about  
841 the converse?

842 **Definition H.2.** A set of real numbers  $L_{ij}$ , one per connection of a given neural architecture, is  
843 self-consistent if and only if there is a unique corresponding set of numbers  $\Lambda_i > 0$  (one per unit)  
844 such that:  $\Lambda_i = 1$  for all visible units and  $L_{ij} = \log \Lambda_i/\Lambda_j$  for every directed connection from a unit  
845  $j$  to a unit  $i$ .

846 **Remark H.3.** This definition depends on the graph of connections, but not on the original values of  
847 the synaptic weights. Every balanced state is associated with a self-consistent set of  $L_{ij}$ , but not  
848 every self-consistent set of  $L_{ij}$  is associated with a balanced state.

849 **Proposition H.4.** A set  $L_{ij}$  associated with a neural architecture is self-consistent if and only if  
850  $\sum_{\pi} L_{ij} = 0$  where  $\pi$  is any directed path connecting an input unit to an output unit or any directed  
851 cycle (for recurrent networks).

852 **Remark H.5.** Thus the constraints associated with being a self-consistent configuration of  $L_{ij}$ 's  
853 are all linear. This linear manifold of constraints depends only on the architecture, i.e., the graph of  
854 connections. The strictly convex function  $R(L_{ij})$  depends on the actual weights  $W$ . Different sets of  
855 weights  $W$  produce different convex functions over the same linear manifold.

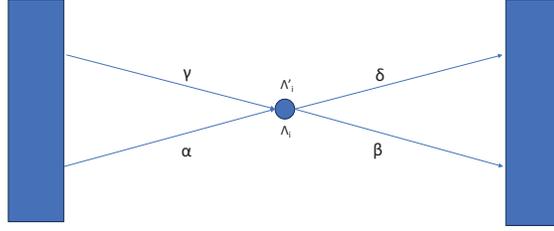


Figure 9: Consider two paths  $\alpha + \beta$  and  $\gamma + \delta$  from the input layer to the output layer going through the same unit  $i$ . Let us assume that the first path assigns a multiplier  $\Lambda_i$  to unit  $i$  and the second path assigns a multiplier  $\Lambda'_i$  to the same unit. By assumption we must have:  $\sum_{\alpha} L_{ij} + \sum_{\beta} L_{ij} = 0$  for the first path, and  $\sum_{\gamma} L_{ij} + \sum_{\delta} L_{ij} = 0$ . But  $\alpha + \delta$  and  $\gamma + \beta$  are also paths from the input layer to the output layer and therefore:  $\sum_{\alpha} L_{ij} + \sum_{\delta} L_{ij} = 0$  and  $\sum_{\gamma} L_{ij} + \sum_{\beta} L_{ij} = 0$ . As a result,  $\sum_{\alpha} L_{ij} = \log \Lambda_i = \sum_{\gamma} L_{ij} = \Lambda'_i$ . Therefore the assignment of the multiplier  $\Lambda_i$  must be consistent across different paths going through unit  $i$ .

856 *Remark H.6.* Note that one could coalesce all the input units and all output units into a single unit,  
 857 in which case a path from an input unit to and output unit becomes also a directed cycle. In this  
 858 representation, the constraints are that the sum of the  $L_{ij}$  must be zero along any directed cycle. In  
 859 general, it is not necessary to write a constraint for every path from input units to output units. It is  
 860 sufficient to select a representative set of paths such that every unit appears in at least one path.

861 *Proof.* If we look at any directed path  $\pi$  from unit  $i$  to unit  $j$ , it is easy to see that we must have:

$$\sum_{\pi} L_{kl} = \log \Lambda_i - \log \Lambda_j \quad (\text{H.7})$$

862 This is illustrated in Figures 8 and 1. Thus along any directed path that connects any input unit to any  
 863 output unit, we must have  $\sum_{\pi} L_{ij} = 0$ . In addition, for recurrent neural networks, if  $\pi$  is a directed  
 864 cycle we must also have:  $\sum_{\pi} L_{ij} = 0$ . Thus in short we only need to add linear constraints of the  
 865 form:  $\sum_{\pi} L_{ij} = 0$ . Any unit is situated on a path from an input unit to an output unit. Along that  
 866 path, it is easy to assign a value  $\Lambda_i$  to each unit by simple propagation starting from the input unit  
 867 which has a multiplier equal to 1. When the propagation terminates in the output unit, it terminates  
 868 consistently because the output unit has a multiplier equal to 1 and, by assumption, the sum of the  
 869 multipliers along the path must be zero. So we can derive scaling values  $\Lambda_i$  from the variables  
 870  $L_{ij}$ . Finally, we need to show that there are no clashes, i.e. that it is not possible for two different  
 871 propagation paths to assign different multiplier values to the same unit  $i$ . The reason for this is  
 872 illustrated in Figure 9.  $\square$

873 We can now complete the proof Theorem H.1. Given a neural network of BiLUs with a set of weights  
 874  $W$ , we can consider the problem of minimizing the regularizer  $R(L_{ij})$  over the self-admissible  
 875 configuration  $L_{ij}$ . For any  $P > 0$ , the  $L_p$  regularizer is strictly convex and the space of self-  
 876 admissible configurations is linear and hence convex. Thus this is a strictly convex optimization  
 877 problem that has a unique solution (Figure 2). Note that the minimization is carried over self-  
 878 consistent configurations, which in general are not associated with balanced states. However, the  
 879 configuration of the weights associated with the optimum set of  $L_{ij}$  (point  $A$  in Figure 2) must be  
 880 balanced. To see this, imagine that one of the BiLU units—unit  $i$  in the network is not balanced. Then  
 881 we can balance it using a multiplier  $\lambda_i^*$  and replace  $\Lambda_i$  by  $\Lambda'_i = \Lambda_i \lambda_i^*$ . It is easy to check that the new  
 882 configuration including  $\Lambda'_i$  is self-consistent. Thus, by balancing unit  $i$ , we are able to reach a new  
 883 self-consistent configuration with a lower value of  $R$  which contradicts the fact that we are at the  
 884 global minimum of the strictly convex optimization problem.

885 We know that the stochastic balancing algorithm always converges to a balanced state. We need to  
886 show that it cannot converge to any other balanced state, and in fact that the global optimum is the  
887 only balanced state. By contradiction, suppose it converges to a different balanced state associated  
888 with the coordinates  $(L_{ij}^B)$  (point  $B$  in Figure 2). Because of the self-consistency, this point is also  
889 associated with a unique set of  $(\Lambda_i^B)$  coordinates. The cost function is continuous and differentiable  
890 in both the  $L_{ij}$ 's and the  $\Lambda_i$ 's coordinates. If we look at the negative gradient of the regularizer, it  
891 is non-zero and therefore it must have at least one non-zero component  $\partial R/\partial \Lambda_i$  along one of the  
892  $\Lambda_i$  coordinates. This implies that by scaling the corresponding unit  $i$  in the network, the regularizer  
893 can be further reduced, and by balancing unit  $i$  the balancing algorithm will reach a new point ( $C$  in  
894 Figure 2) with lower regularizer cost. This contradicts the assumption that  $B$  was associated with a  
895 balanced state. Thus, given an initial set of weights  $W$ , the stochastic balancing algorithm must  
896 always converge to the same and unique optimal balanced state  $W^*$  associated with the self-consistent  
897 point  $A$ . A particular stochastic schedule corresponds to a random path within the linear manifold  
898 from the origin (at time zero all the multipliers are equal to 1, and therefore  $M_{ij} = 1$  and  $L_{ij} = 0$ )  
899 for any  $i$  and any  $j$  to the unique optimum point  $A$ .  $\square$

900 *Remark H.7.* It should be clear from the proof that the same result holds also for any deterministic  
901 balancing schedule, as well as for tied and non-tied subset balancing, e.g., for layer-wise balancing  
902 and tied layer-wise balancing. In the Appendix, we provide an analytical solution for the case of tied  
903 layer-wise balancing in a layered feed-forward network.

904 *Remark H.8.* It should be clear from the proof that the same convergence to the unique global  
905 optimum is observed if each neuron, when stochastically visited, is favorably scaled rather than  
906 balanced, i.e., it is scaled with a factor that reduces  $R$  but not necessarily minimizes  $R$ . Stochastic  
907 balancing can also be viewed as a form of EM algorithm where the E and M steps can be taken fully  
908 or partially.

## 909 I Universal Approximation Properties of BiLU Neurons

910 Here we show that any continuous real-valued function defined over a compact set of the Euclidean  
911 space can be approximated to any degree of precision by a network of BiLU neurons with a single  
912 hidden layer. As in the case of the similar proof given in Baldi [2021] using linear threshold gates in  
913 the hidden layer, it is enough to prove the theorem for a continuous function  $f: [0, 1] \rightarrow \mathbb{R}$ .

914 **Theorem I.1.** (*Universal Approximation Properties of BiLU Neurons*) *Let  $f$  be any continuous*  
915 *function from  $[0, 1]$  to  $\mathbb{R}$  and  $\epsilon > 0$ . Let  $g_\lambda$  be the ReLU activation function with slope  $\lambda \in \mathbb{R}$ s. Then*  
916 *there exists a feedforward network with a single hidden layer of neurons with ReLU activations of the*  
917 *form  $g_\lambda$  and a single output linear neuron, i.e., with BiLU activation equal to the identity function,*  
918 *capable of approximating  $f$  everywhere within  $\epsilon$  (sup norm).*

919 *Proof.* To be clear,  $g_\lambda(x) = 0$  for  $x < 0$  and  $g_\lambda(x) = \lambda x$  for  $0 \leq x$ . Since  $f$  is continuous over a  
920 compact set, it is uniformly continuous. Thus there exists  $\alpha > 0$  such that for any  $x_1$  and  $x_2$  in the  
921  $[0, 1]$  interval:

$$|x_2 - x_1| < \alpha \implies |f(x_2) - f(x_1)| < \epsilon \quad (\text{I.1})$$

922 Let  $N$  be an integer such that  $1 < N\alpha$ , and let us slice the interval  $[0, 1]$  into  $N$  consecutive slices  
923 of width  $h = 1/N$ , so that within each slice the function  $f$  cannot jump by more than  $\epsilon$ . Let us  
924 connect the input unit to all the hidden units with a weight equal to 1. Let us have  $N$  hidden units  
925 numbered  $1, \dots, N$  with biases equal to  $0, 1/N, 2/N, \dots, N_1/N$  respectively and activation function  
926 of the form  $g_{\lambda_k}$ . It is essential that different units be allowed to have different slopes  $\lambda_k$ . The input  
927 unit is connected to all the hidden units and all the weights on these connections are equal to 1. Thus  
928 when  $x$  is in the  $k$ -th slice,  $(k-1)/N \leq x < k/N$ , all the units from  $k+1$  to  $N$  have an output  
929 equal to 0, and all the units from 1 to  $k$  have an output determined by the corresponding slopes. All  
930 the hidden units are connected to the output unit with weights  $\beta_1, \dots, \beta_N$ , and  $\beta_0$  is the bias of the  
931 output unit. We want the output unit to be linear. In order for the  $\epsilon$  approximation to be satisfied,  
932 it is sufficient if in the  $(k-1)/N \leq x < k/N$  interval, the output is equal to the line joining the  
933 point  $f((k-1)/N)$  to the point  $f(k/N)$ . In other words, if  $x \in [(k-1)/N, k/N)$ , then we want  
934 the output of the network to be:

$$\beta_0 + \sum_{i=1}^k \beta_i \lambda_i (x - (i-1)h) = f\left(\frac{k-1}{N}\right) + \frac{f\left(\frac{k}{N}\right) - f\left(\frac{k-1}{N}\right)}{h} (x - (k-1)h) \quad (\text{I.2})$$

935 By equating the y-intercept and slope of the lines on the left-hand side and the right-hand side of  
 936 Equation I.2, we can solve for the weights  $\beta$ 's and the slopes  $\lambda$ 's.  $\square$

937 As in the case of the similar proof using linear threshold functions in the hidden layer (see Baldi  
 938 [2021]), this proof can easily be adapted to continuous functions defined over a compact set of  $\mathbb{R}^n$ ,  
 939 even with a finite number of finite discontinuities, and into  $\mathbb{R}^m$ .

## 940 J Analytical Solution for the Unique Global Balanced State

941 Here we directly prove the convergence of stochastic balancing to a unique final balanced state, and  
 942 derive the equations for the balanced state, in the special case of tied layer balancing (as opposed to  
 943 single neuron balancing). The Proof and the resulting equations are also valid for stochastic balancing  
 944 (one neuron at a time) in a layered architecture comprising a single neuron per layer. Let us call tied  
 945 layer scaling the operation by which all the incoming weights to a given layer of BiLU neurons are  
 946 multiplied by  $\lambda > 0$  and all the outgoing weights of the layer are multiplied by  $1/\lambda$ , again leaving the  
 947 training error unchanged. Let us call layer balancing the particular scaling operation corresponding  
 948 to the value of  $\lambda$  that minimizes the contribution of the layer to the  $L_2$  (or any other  $L_p$ ) regularizer  
 949 value. This optimal value of  $\lambda^*$  results in layer-wise balance equations: the sum of the squares of all  
 950 the incoming weights of the layer must be equal to the sum of the squares of all the outgoing weights  
 951 of the layer in the  $L_2$  case, and similarly in all  $L^p$  cases.

952 **Theorem J.1.** *Assume that tied layer balancing is applied iteratively and stochastically to the layers*  
 953 *of a layered feedforward network of BiLU neurons. As long as all the layers are visited periodically,*  
 954 *this procedure will always converge to the same unique set of weights, which will satisfy the layer-*  
 955 *balance equations at all layers, irrespective of the details of the schedule. Furthermore, the balance*  
 956 *state can be solved analytically.*

957 *Proof.* Every time a layer balancing operation is applied, the training error remains the same, and the  
 958  $L_2$  (or any other  $L_p$ ) regularization error decreases or stays the same. Since the regularization error  
 959 is always positive, it must converge to a certain value. Using the same arguments as in the proof of  
 960 Theorem H.1, the weights must also converge to a stable configuration, and since the configuration  
 961 is stable all its layers must satisfy the layer-wise balance equation. The key remaining question is  
 962 why is this configuration unique and can we solve it analytically? Let  $A_1, A_2, \dots, A_N$  denote the  
 963 matrices of connections between the layers of the network. Let  $\Lambda_1, \Lambda_2, \dots, \Lambda_{N-1}$  be  $N-1$  strictly  
 964 positive multipliers, representing the limits of the products of the corresponding  $\lambda_i^*$  associated with  
 965 each balancing step at layer  $i$ , as in the proof of Theorem H.1. In this notation, layer 0 is the input  
 966 layer and layer  $N$  is the output layer (with  $\Lambda_0 = 1$  and  $\Lambda_N = 1$ ).

967 After converging, each matrix  $A_i$  becomes the matrix  $\Lambda_i/\Lambda_{i-1}A_i = M_iA_i$  for  $i = 1 \dots N$ , with  
 968  $M_i = \lambda_i/\Lambda_{i-1}$ . The multipliers  $M_i$  must minimize the regularizer while satisfying  $M_1 \dots M_N = 1$   
 969 to ensure that the training error remains unchanged. In other words, to find the values of the  $M_i$ 's we  
 970 must minimize the Lagrangian:

$$\mathcal{L}(M_1, \dots, M_N) = \sum_{i=1}^N \|M_i A_i\|^2 + \mu \left(1 - \prod_{i=1}^N M_i\right) \quad (\text{J.1})$$

971 written for the  $L^2$  case in terms of the Frobenius norm, but the analysis is similar in the general  $L_p$   
 972 case. From this, we get the critical equations:

$$\frac{\partial \mathcal{L}}{\partial M_i} = 2M_i \|A_i\|^2 - \mu M_1 \dots M_{i-1} M_{i+1} \dots M_N = 0 \quad \text{for } i = 1, \dots, N \quad \text{and} \quad \prod_{i=1}^N M_i = 1 \quad (\text{J.2})$$

973 As a result, for every  $i$ :

$$2M_i||A_i||^2 - \frac{\mu}{M_i} = 0 \quad \text{or} \quad \mu = 2M_i^2||A_i||^2 \quad (\text{J.3})$$

974 Thus each  $M_i > 0$  can be expressed in a unique way as a function of the Lagrangian multiplier  $\mu$  as:  
 975  $M_i = (\mu/2||A_i||^2)^{1/2}$ . By writing again that the product of the  $M_i$  is equal to 1, we finally get:

$$\mu^N = 2^N \prod_{i=1}^N ||A_i||^2 \quad \text{or} \quad \mu = 2 \prod_{i=1}^N ||A_i||^{2/N} \quad (\text{J.4})$$

976 Thus we can solve for  $M_i$ :

$$M_i = \frac{\mu}{2||A_i||^2} = \frac{\prod_{i=1}^N ||A_i||^{2/N}}{||A_i||^2} \quad \text{for } i = 1, \dots, N \quad (\text{J.5})$$

977 Thus, in short, we obtain a unique closed-form expression for each  $M_i$ . From there, we infer the  
 978 unique and final state of the weights, where  $A_i^* = M_i A_i = \Lambda_i A_i / \Lambda_{l-1}$ . Note that each  $M_i$  depends  
 979 on all the other  $M_j$ 's, again showcasing how the local balancing algorithm leads to a unique global  
 980 solution.  $\square$

## 981 **K Computer Resources**

982 The simulations we have described do not require major computing resources. They were all  
 983 performed using Google Colab and the NVIDIA TESLA T4 GPU that it provides.

## 984 **L Code Availability**

985 The code for reproducing the simulation results is available under the Apache 2.0 license at:  
 986 <https://anonymous.4open.science/r/a-theory-of-neural-synaptic-balance-00C1>

987 **References**

- 988 P. Baldi. *Deep Learning in Science*. Cambridge University Press, Cambridge, UK, 2021.
- 989 Li Deng. The mnist database of handwritten digit images for machine learning research. *IEEE Signal*  
990 *Processing Magazine*, 29(6):141–142, 2012.
- 991 Simon S Du, Wei Hu, and Jason D Lee. Algorithmic regularization in learning deep homogeneous  
992 models: Layers are automatically balanced. *Advances in Neural Information Processing Systems*,  
993 31, 2018.
- 994 Rachel E Field, James A D’amour, Robin Tremblay, Christoph Miehl, Bernardo Rudy, Julijana  
995 Gjorgjieva, and Robert C Froemke. Heterosynaptic plasticity determines the set point for cortical  
996 excitatory-inhibitory balance. *Neuron*, 106(5):842–854, 2020.
- 997 Robert C Froemke. Plasticity of cortical excitatory-inhibitory balance. *Annual review of neuroscience*,  
998 38:195–219, 2015.
- 999 Oliver D Howes and Ekaterina Shatalina. Integrating the neurodevelopmental and dopamine hypothe-  
1000 ses of schizophrenia and the role of cortical excitation-inhibition balance. *Biological psychiatry*,  
1001 2022.
- 1002 Dmitry Ivanov, Aleksandr Chezhegov, Mikhail Kiselev, Andrey Grunin, and Denis Larionov. Neuro-  
1003 morphic artificial intelligence systems. *Frontiers in Neuroscience*, 16:1513, 2022.
- 1004 Yu Ji, YouHui Zhang, ShuangChen Li, Ping Chi, CiHang Jiang, Peng Qu, Yuan Xie, and WenGuang  
1005 Chen. Neutrams: Neural network transformation and co-design under neuromorphic hardware  
1006 constraints. In *2016 49th Annual IEEE/ACM International Symposium on Microarchitecture*  
1007 *(MICRO)*, pages 1–13. IEEE, 2016.
- 1008 Dongshin Kim and Jang-Sik Lee. Neurotransmitter-induced excitatory and inhibitory functions in  
1009 artificial synapses. *Advanced Functional Materials*, 32(21):2200497, 2022.
- 1010 Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. 2009.
- 1011 Faming Liang and Wing Hung Wong. Evolutionary monte carlo: Applications to cp model sampling  
1012 and change point problem. *STATISTICA SINICA*, 10:317–342, 2000.
- 1013 Behnam Neyshabur, Ryota Tomioka, Ruslan Salakhutdinov, and Nathan Srebro. Data-dependent path  
1014 normalization in neural networks. *arXiv preprint arXiv:1511.06747*, 2015.
- 1015 Bodo Rueckauer, Iulia-Alexandra Lungu, Yuhuang Hu, Michael Pfeiffer, and Shih-Chii Liu. Con-  
1016 version of continuous-valued deep networks to efficient event-driven networks for image classification.  
1017 *Frontiers in neuroscience*, 11:294078, 2017.
- 1018 Farshad Shirani and Hannah Choi. On the physiological and structural contributors to the dynamic  
1019 balance of excitation and inhibition in local cortical networks. *bioRxiv*, pages 2023–01, 2023.
- 1020 Martino Sorbaro, Qian Liu, Massimo Bortone, and Sadique Sheik. Optimizing the energy consump-  
1021 tion of spiking neural networks for neuromorphic applications. *Frontiers in neuroscience*, 14:662,  
1022 2020.
- 1023 Christopher H Stock, Sarah E Harvey, Samuel A Ocko, and Surya Ganguli. Synaptic balancing: A  
1024 biologically plausible local learning rule that provably increases neural network noise robustness  
1025 without sacrificing task performance. *PLOS Computational Biology*, 18(9):e1010418, 2022.
- 1026 A. Tavakoli, F. Agostinelli, and P. Baldi. SPLASH: Learnable activation functions for improving  
1027 accuracy and adversarial robustness. *Neural Networks*, 140:1–12, 2021. Also: arXiv:2006.08947.

## 1028 **NeurIPS Paper Checklist**

### 1029 **1. Claims**

1030 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1031 paper's contributions and scope?

1032 Answer: [Yes]

1033 Justification: We have included all the main points of the paper in the abstract.

1034 Guidelines:

- 1035 • The answer NA means that the abstract and introduction do not include the claims  
1036 made in the paper.
- 1037 • The abstract and/or introduction should clearly state the claims made, including the  
1038 contributions made in the paper and important assumptions and limitations. A No or  
1039 NA answer to this question will not be perceived well by the reviewers.
- 1040 • The claims made should match theoretical and experimental results, and reflect how  
1041 much the results can be expected to generalize to other settings.
- 1042 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1043 are not attained by the paper.

### 1044 **2. Limitations**

1045 Question: Does the paper discuss the limitations of the work performed by the authors?

1046 Answer: [Yes]

1047 Justification: The majority of our results are theorems backed up by mathematical proofs.  
1048 We discuss at length that balancing improves the value of the regularizer only (it leaves the  
1049 value of the data-dependent component of the error unchanged). We also mention that  
1050 while it would be interesting to study any kind of balance in biological neural networks,  
1051 current technological limitations do not allow recording all the incoming and outgoing  
1052 synaptic strengths of a neuron.

1053 Guidelines:

- 1054 • The answer NA means that the paper has no limitation while the answer No means that  
1055 the paper has limitations, but those are not discussed in the paper.
- 1056 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1057 • The paper should point out any strong assumptions and how robust the results are to  
1058 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1059 model well-specification, asymptotic approximations only holding locally). The authors  
1060 should reflect on how these assumptions might be violated in practice and what the  
1061 implications would be.
- 1062 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1063 only tested on a few datasets or with a few runs. In general, empirical results often  
1064 depend on implicit assumptions, which should be articulated.
- 1065 • The authors should reflect on the factors that influence the performance of the approach.  
1066 For example, a facial recognition algorithm may perform poorly when image resolution  
1067 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1068 used reliably to provide closed captions for online lectures because it fails to handle  
1069 technical jargon.
- 1070 • The authors should discuss the computational efficiency of the proposed algorithms  
1071 and how they scale with dataset size.
- 1072 • If applicable, the authors should discuss possible limitations of their approach to  
1073 address problems of privacy and fairness.
- 1074 • While the authors might fear that complete honesty about limitations might be used by  
1075 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1076 limitations that aren't acknowledged in the paper. The authors should use their best  
1077 judgment and recognize that individual actions in favor of transparency play an impor-  
1078 tant role in developing norms that preserve the integrity of the community. Reviewers  
1079 will be specifically instructed to not penalize honesty concerning limitations.

### 1080 **3. Theory Assumptions and Proofs**

1081 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1082 a complete (and correct) proof?

1083 Answer: [Yes]

1084 Justification: All the theorems and propositions have clear assumptions and all the proofs  
1085 are complete and have been checked carefully multiple times. Details of some of the proofs  
1086 are provided in the Appendix.

1087 Guidelines:

- 1088 • The answer NA means that the paper does not include theoretical results.
- 1089 • All the theorems, formulas, and proofs in the paper should be numbered and cross-  
1090 referenced.
- 1091 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 1092 • The proofs can either appear in the main paper or the supplemental material, but if  
1093 they appear in the supplemental material, the authors are encouraged to provide a short  
1094 proof sketch to provide intuition.
- 1095 • Inversely, any informal proof provided in the core of the paper should be complemented  
1096 by formal proofs provided in appendix or supplemental material.
- 1097 • Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 1098 4. Experimental Result Reproducibility

1099 Question: Does the paper fully disclose all the information needed to reproduce the main ex-  
1100 perimental results of the paper to the extent that it affects the main claims and/or conclusions  
1101 of the paper (regardless of whether the code and data are provided or not)?

1102 Answer: [Yes]

1103 Justification: We have provided all the explanations necessary for reproducing the exper-  
1104 imental results in the technical appendix and also provided the code for reproducing our  
1105 experimental results.

1106 Guidelines:

- 1107 • The answer NA means that the paper does not include experiments.
- 1108 • If the paper includes experiments, a No answer to this question will not be perceived  
1109 well by the reviewers: Making the paper reproducible is important, regardless of  
1110 whether the code and data are provided or not.
- 1111 • If the contribution is a dataset and/or model, the authors should describe the steps taken  
1112 to make their results reproducible or verifiable.
- 1113 • Depending on the contribution, reproducibility can be accomplished in various ways.  
1114 For example, if the contribution is a novel architecture, describing the architecture fully  
1115 might suffice, or if the contribution is a specific model and empirical evaluation, it may  
1116 be necessary to either make it possible for others to replicate the model with the same  
1117 dataset, or provide access to the model. In general, releasing code and data is often  
1118 one good way to accomplish this, but reproducibility can also be provided via detailed  
1119 instructions for how to replicate the results, access to a hosted model (e.g., in the case  
1120 of a large language model), releasing of a model checkpoint, or other means that are  
1121 appropriate to the research performed.
- 1122 • While NeurIPS does not require releasing code, the conference does require all submis-  
1123 sions to provide some reasonable avenue for reproducibility, which may depend on the  
1124 nature of the contribution. For example
  - 1125 (a) If the contribution is primarily a new algorithm, the paper should make it clear how  
1126 to reproduce that algorithm.
  - 1127 (b) If the contribution is primarily a new model architecture, the paper should describe  
1128 the architecture clearly and fully.
  - 1129 (c) If the contribution is a new model (e.g., a large language model), then there should  
1130 either be a way to access this model for reproducing the results or a way to reproduce  
1131 the model (e.g., with an open-source dataset or instructions for how to construct  
1132 the dataset).

1133 (d) We recognize that reproducibility may be tricky in some cases, in which case  
1134 authors are welcome to describe the particular way they provide for reproducibility.  
1135 In the case of closed-source models, it may be that access to the model is limited in  
1136 some way (e.g., to registered users), but it should be possible for other researchers  
1137 to have some path to reproducing or verifying the results.

## 1138 5. Open access to data and code

1139 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
1140 tions to faithfully reproduce the main experimental results, as described in supplemental  
1141 material?

1142 Answer: [Yes]

1143 Justification: We have provided an anonymous link to our code which is available in the  
1144 appendix and also uploaded our code as supplementary material.

1145 Guidelines:

- 1146 • The answer NA means that paper does not include experiments requiring code.
- 1147 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
1148 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1149 • While we encourage the release of code and data, we understand that this might not be  
1150 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
1151 including code, unless this is central to the contribution (e.g., for a new open-source  
1152 benchmark).
- 1153 • The instructions should contain the exact command and environment needed to run to  
1154 reproduce the results. See the NeurIPS code and data submission guidelines ([https:  
1155 //nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 1156 • The authors should provide instructions on data access and preparation, including how  
1157 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 1158 • The authors should provide scripts to reproduce all experimental results for the new  
1159 proposed method and baselines. If only a subset of experiments are reproducible, they  
1160 should state which ones are omitted from the script and why.
- 1161 • At submission time, to preserve anonymity, the authors should release anonymized  
1162 versions (if applicable).
- 1163 • Providing as much information as possible in supplemental material (appended to the  
1164 paper) is recommended, but including URLs to data and code is permitted.

## 1165 6. Experimental Setting/Details

1166 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
1167 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
1168 results?

1169 Answer: [Yes]

1170 Justification: We have provided the required details in the appendix.

1171 Guidelines:

- 1172 • The answer NA means that the paper does not include experiments.
- 1173 • The experimental setting should be presented in the core of the paper to a level of detail  
1174 that is necessary to appreciate the results and make sense of them.
- 1175 • The full details can be provided either with the code, in appendix, or as supplemental  
1176 material.

## 1177 7. Experiment Statistical Significance

1178 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
1179 information about the statistical significance of the experiments?

1180 Answer: [Yes]

1181 Justification: Error bars are included in all images.

1182 Guidelines:

- 1183 • The answer NA means that the paper does not include experiments.

- 1184 • The authors should answer "Yes" if the results are accompanied by error bars, confi-  
1185 dence intervals, or statistical significance tests, at least for the experiments that support  
1186 the main claims of the paper.
- 1187 • The factors of variability that the error bars are capturing should be clearly stated (for  
1188 example, train/test split, initialization, random drawing of some parameter, or overall  
1189 run with given experimental conditions).
- 1190 • The method for calculating the error bars should be explained (closed form formula,  
1191 call to a library function, bootstrap, etc.)
- 1192 • The assumptions made should be given (e.g., Normally distributed errors).
- 1193 • It should be clear whether the error bar is the standard deviation or the standard error  
1194 of the mean.
- 1195 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1196 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1197 of Normality of errors is not verified.
- 1198 • For asymmetric distributions, the authors should be careful not to show in tables or  
1199 figures symmetric error bars that would yield results that are out of range (e.g. negative  
1200 error rates).
- 1201 • If error bars are reported in tables or plots, The authors should explain in the text how  
1202 they were calculated and reference the corresponding figures or tables in the text.

## 1203 8. Experiments Compute Resources

1204 Question: For each experiment, does the paper provide sufficient information on the com-  
1205 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1206 the experiments?

1207 Answer: [Yes]

1208 Justification: We have provided this information in the computer resources section in the  
1209 appendix.

1210 Guidelines:

- 1211 • The answer NA means that the paper does not include experiments.
- 1212 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1213 or cloud provider, including relevant memory and storage.
- 1214 • The paper should provide the amount of compute required for each of the individual  
1215 experimental runs as well as estimate the total compute.
- 1216 • The paper should disclose whether the full research project required more compute  
1217 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1218 didn't make it into the paper).

## 1219 9. Code Of Ethics

1220 Question: Does the research conducted in the paper conform, in every respect, with the  
1221 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1222 Answer: [Yes]

1223 Justification: The research conducted in our paper conforms, in every respect, with the  
1224 NeurIPS Code of Ethics.

1225 Guidelines:

- 1226 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1227 • If the authors answer No, they should explain the special circumstances that require a  
1228 deviation from the Code of Ethics.
- 1229 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
1230 eration due to laws or regulations in their jurisdiction).

## 1231 10. Broader Impacts

1232 Question: Does the paper discuss both potential positive societal impacts and negative  
1233 societal impacts of the work performed?

1234 Answer: [NA]

1235 Justification: Our paper has no conceivable direct societal impact.

1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: Our paper poses no such risks.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The only assets that we have use are the MNIST and CIFAR-10 datasets and we have cited these datasets in the paper properly.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.

- 1289
- 1290
- 1291
- 1292
- 1293
- 1294
- 1295
- 1296
- 1297
- 1298
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
  - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
  - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
  - If this information is not available online, the authors are encouraged to reach out to the asset's creators.

1299 **13. New Assets**

1300 Question: Are new assets introduced in the paper well documented and is the documentation  
1301 provided alongside the assets?

1302 Answer: [NA]

1303 Justification: Our paper does not introduce new assets.

1304 Guidelines:

- 1305
- 1306
- 1307
- 1308
- 1309
- 1310
- 1311
- 1312
- The answer NA means that the paper does not release new assets.
  - Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
  - The paper should discuss whether and how consent was obtained from people whose asset is used.
  - At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

1313 **14. Crowdsourcing and Research with Human Subjects**

1314 Question: For crowdsourcing experiments and research with human subjects, does the paper  
1315 include the full text of instructions given to participants and screenshots, if applicable, as  
1316 well as details about compensation (if any)?

1317 Answer: [NA]

1318 Justification: Our research does not involve human subjects or crowdsourcing.

1319 Guidelines:

- 1320
- 1321
- 1322
- 1323
- 1324
- 1325
- 1326
- 1327
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
  - Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
  - According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

1328 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human  
1329 Subjects**

1330 Question: Does the paper describe potential risks incurred by study participants, whether  
1331 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
1332 approvals (or an equivalent approval/review based on the requirements of your country or  
1333 institution) were obtained?

1334 Answer: [NA]

1335 Justification: Our research does not involve any human subjects.

1336 Guidelines:

- 1337
- 1338
- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.