



# A cookbook for hardware-friendly implicit learning on static data

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 The following aims to be a pragmatic introduction to hardware-friendly learning of  
2 *implicit models*, which encompass a broad class of models from feedforward nets  
3 to physical systems, taking static data as inputs. Starting from first principles, we  
4 present a *minimal* hierarchy of independent concepts to circumvent some problems  
5 inherent to the hardware implementation of standard differentiation. This way,  
6 we avoid entangling essential ingredients with arbitrary design choices by naively  
7 listing existing algorithms and instead propose the draft of a “cookbook” to help  
8 the exploration of many possible combinations of these independent mechanisms.

## 9 1 Problem statement

10 **Learning at equilibrium.** Given an input  $x$ , we want to find a set of model parameters  $\theta$  which  
11 minimizes a given objective  $\mathcal{O}$  defined over the model (hidden and output) variables  $s(x, \theta)$ , such  
12 that the model variables abides by some constraint  $\mathcal{C}$ . Namely:

$$\mathcal{P}_1 : \min_{\theta} J(x, \theta) := \mathcal{O}(s(x, \theta), \theta) \quad \text{s.t.} \quad \mathcal{C}(x, s(x, \theta), \theta) = 0. \quad (1)$$

13 Note that  $s$  may implicitly contain several layers, e.g.  $s = (s^1{}^\top, s^2{}^\top, \dots, s^N{}^\top)^\top$ . Classically,  $\mathcal{O}$   
14 is some *cost function*  $\ell$  measuring the discrepancy between the model prediction and some ground-  
15 truth data and  $\mathcal{C}$  the “physical laws” that governs the substrate sustaining the model prediction *at*  
16 *equilibrium* – note that this implicit formalism encompasses both feedforward models (60), deep  
17 equilibrium models (4), as well as resistive networks governed by Kirchoff’s laws (39; 86; 77).  
18 The problem defined in Eq. (1) is traditionally solved via first-order optimization by estimating the  
19 gradient over a minibatch  $\mathcal{B}$  (93):  $\bar{g}(\theta) := \mathbb{E}_{x \sim \mathcal{D}}[d_{\theta} J(x, \theta)]$ . Therefore, solving  $\mathcal{P}_1$  boils down to  
20 *how to compute* these gradients. We herein call a “circuit” the abstract physical system that may  
21 realize a given equilibrium condition.

22 **The Lagrangian method.** One way to solve  $\mathcal{P}_1$  is by writing the associated *Lagrangian* functional  
23  $\mathcal{L}_1(s, \lambda, \theta) := \mathcal{O}(s, \theta) + \lambda^\top \cdot \mathcal{C}(s, \theta)$  and solving for the *Karush-Kuhn-Tucker* (KKT) conditions  
24 (81; 12). Loosely speaking, the *primal* feasibility condition  $\partial_{\lambda} \mathcal{L}_1 = 0$  and *dual* feasibility condition  
25  $\partial_s \mathcal{L}_1 = 0$  generalize the notions of “forward pass” and “backward pass” respectively, and can be  
26 viewed as two circuits determining  $s$  and  $\lambda$  whose equilibrium satisfy:

$$\begin{cases} \mathcal{C}(s, \theta) = 0 \\ \partial_s \mathcal{C}(s, \theta)^\top \cdot \lambda + \nabla_s \mathcal{O}(s, \theta) = 0 \end{cases} \quad (2)$$

27 with the resulting gradient estimate reading as  $g(\theta) = \nabla_{\theta} \mathcal{O}(s, \theta) + \partial_{\theta} \mathcal{C}(s, \theta)^\top \cdot \lambda$ . We briefly  
28 mention that this can also be derived using the *Implicit Function Theorem* (IFT) (93). Here as well,

29 note that *feedforward nets trained by backprop* are a special case of Eq. (2). Yet, it appears from  
 30 Eq. (2) that naively solving for  $s$  and  $\lambda$  comes with some challenges from the (high-level) hardware  
 31 viewpoint: i) **the transport problem**:  $s$  and  $\theta$  need to be transported from the inference circuit to  
 32 the error circuit; ii) **the memory problem**:  $s$  needs to be stored; iii) **analytical derivatives**:  $\partial_s \mathcal{C}$   
 33 potentially contains analytical derivatives which may have to be computed with high precision; iv)  
 34 **forward locking**: in a hierarchical model, the circuit computing  $s^k$  becomes idle when computing  
 35  $s^{k+1}$  (Eq. (3)); v) **backward locking**: symmetrically,  $\lambda_k$  cannot be solved before  $\lambda_{k+1}$  (Eq. (4)); vi)  
 36 **forward-backward synchrony**:  $\lambda$  is computed *after*  $s$ .

## 37 2 Algorithm design & general methods

### 38 2.1 Crafting the constraints

39 **Hierarchical constraints.** As the variable  $s$  may subsume a hierarchy of layers, one may want to  
 40 split the original optimization problem into a hierarchy of  $K$  problems, given  $s := (s^{1^\top}, \dots, s^{K^\top})^\top$   
 41 (15; 47; 29; 92; 69):

$$\mathcal{P}_2 : \min_{\theta} J(x, \theta) := \mathcal{O}(s, \theta) \quad \text{s.t.} \quad \mathcal{C}_1(s^1, s^0 := x, \theta^1) = 0, \dots, \mathcal{C}_K(s^K, s^{K-1}, \theta^K) = 0 \quad (3)$$

42 Here, the model is split into a hierarchy of  $K$  subcircuits (or “blocks”) with parameters  $\theta^k$ , state  $s^k$ ,  
 43 which may comprise one *or multiple layers* ( $K \leq L$ ), subject to the influence of the previous circuit  
 44 through  $s^{k-1}$ . Note that for a given set of constraints,  $\mathcal{P}_1$  is *generally not equivalent* to  $\mathcal{P}_2$  – see § 3.6  
 45 for an example. In this case, solving for the KKT conditions of  $\mathcal{P}_2$  yields:

$$\partial_{s^k} \mathcal{C}_k(s^k, s^{k-1}, \theta^k)^\top \cdot \lambda^k + \nabla_{s^k} \mathcal{O}(s, \theta) + \partial_{s^k} \mathcal{C}_{k+1}(s^{k+1}, s^k, \theta^{k+1})^\top \cdot \lambda^{k+1} = 0. \quad (4)$$

46 Note that  $\mathcal{O}(s, \theta) = \ell(s^K, y)$  corresponds to the classical “end-to-end” supervised learning setting in  
 47 which case the second term of Eq. (4) vanishes except for the last block  $s^K$ .

48 **Relaxed constraints.** When  $\mathcal{C} = \nabla_s \mathcal{K}$  (see Section 3),  $\mathcal{P}_1$  can be relaxed through an optimal value  
 49 reformulation with a *fixed* Lagrangian multiplier  $\beta^{-1}$  as (73; 92; 35):

$$\mathcal{P}_3 : \min_{\theta, s} \mathcal{O}(s, \theta) \quad \text{s.t.} \quad \mathcal{K}(s, \theta) \leq \min_{s'} \mathcal{K}(s', \theta), \quad \mathcal{L}_3(s, \theta, \beta) := \mathcal{O}(s, \theta) + \frac{1}{\beta} \left( \mathcal{K}(s, \theta) - \min_{s'} \mathcal{K}(s', \theta) \right) \quad (5)$$

50 The resulting Lagrangian  $\mathcal{L}_3$ , sometimes called “surrogate”, is intimately tied to energy-based learning  
 51 (80) (Section 3) and is especially convenient to express quantities as *finite* differences in  $\beta$  which  
 52 would otherwise appear as *exact* derivatives in  $\beta$  starting from  $\mathcal{P}_1$  (which amounts to send  $\beta \rightarrow 0$ ).  
 53 Importantly, this relaxation could also well be applied *at each level of the hierarchy* of  $\mathcal{P}_2$  (35).

### 54 2.2 Picking an algorithm to estimate $\lambda$ ’s: standard approaches

55 **Implicit differentiation.** In spite of the aforementioned problems inherent to Eq. (2), one may  
 56 still directly solve for  $\lambda$  in combination with other tricks (see section 2.3). For feedforward models,  
 57  $\partial_s \mathcal{C}^\top$  may be *explicitly* inverted with  $\lambda = -\partial_s \mathcal{C}^{-\top} \cdot \mathcal{O}$  reducing to *backprop*. In other cases,  $\lambda$   
 58 may be computed by *implicit differentiation* (41; 10; 28), which comes in many different flavors  
 59 depending on the constraints at hand (13). Note that these techniques can be slightly adjusted to  
 60 accommodate some constraints. For instance, to avoid weight transport between the inference and  
 61 error circuits, we can equip the error circuit with its own “feedback” parameters  $\omega$  such that the error  
 62 signal satisfies:  $\partial_s \mathcal{C}(s, \omega)^\top \cdot \lambda + \nabla_s \mathcal{O}(s, \omega) = 0$  (48; 70; 1; 45). While randomly sampled  $\omega$  are  
 63 sufficient for shallow architectures and simple tasks, some extra alignment mechanisms are needed to  
 64 have  $\omega$  roughly approximate  $\theta$  *without explicitly transporting it* (65; 6; 90; 1; 44).

65 **Forward-only learning with zeroth order optimization.** One technique to solve the memory and  
 66 transport problems, which has recently regained some popularity, is *zeroth order optimization* (ZO)  
 67 (49). ZO techniques estimate a projection of the gradient  $g(\theta)$  along some direction  $u$  in the weight  
 68 space by performing multiple forward passes:

$$u^\top \cdot g(\theta) = d_\epsilon (J(\theta + \epsilon u))|_{\epsilon=0} \approx \frac{1}{2\epsilon} (\mathcal{O}(s(\theta + \epsilon u), \theta + \epsilon u) - \mathcal{O}(s(\theta - \epsilon u), \theta - \epsilon u)), \quad (6)$$

69 An unbiased estimate of the gradient can be obtained by averaging multiple such derivatives:  $g(\theta) =$   
70  $\mathbb{E}_{u \sim \mathcal{N}(0, \sigma^2)} [d_\epsilon (J(\theta + \epsilon u))|_{\epsilon=0} \cdot u]$  (85; 7). However, the variance of this gradient estimate scales  
71 cubically with the number of model parameters (74), thereby restricting the applicability of ZO to the  
72 realm of models which are small enough (84; 33; 7; 24) or behaving as such, i.e. *pre-trained models*  
73 (51). One way to mitigate this problem is to perturb neurons instead of synapses (23; 61; 14; 40),  
74 namely computing projections of the error signal as:  $v^\top \cdot \lambda = d_\epsilon \mathcal{O}(s^\epsilon, \theta)|_{\epsilon=0}$  with  $s^\epsilon$  implicitly  
75 determined through  $\mathcal{C}(s^\epsilon, \theta) + \epsilon v = 0$ . Similarly, an unbiased gradient estimate is obtained as  
76  $g(\theta) = \nabla_\theta \mathcal{O} + \partial_\theta \mathcal{C}(s, \theta)^\top \cdot \mathbb{E}_{v \sim \mathcal{N}(0, \sigma^2)} [d_\epsilon \mathcal{O}(s^\epsilon, \theta)|_{\epsilon=0} \cdot v]$ . Going beyond, one can teach small  
77 auxiliary networks synthesizing such directions to output “good” weight directions  $u$  instead of  
78 randomly sampling them (24), reducing the variance of the gradient estimate at the cost of increasing  
79 the bias (75; 84).

## 80 2.3 Other tricks which independently apply

81 **Greedy learning & loss design.** It appears from Eq. (4) that in general, an error signal  $\lambda^{k+1}$  must  
82 be passed backwards to “unlock” the computation of  $\lambda^k$ . In order to parallelize learning across blocks  
83 entirely, a heuristic consists in shutting off the top-down error signal ( $\lambda^{k+1} = 0$ ) and recreate an error  
84 signal through a *locally-defined* (supervised (8; 71; 9; 88; 27) or self-supervised (50; 91; 31; 83))  
85 loss, which amounts to choosing  $\mathcal{O} = \mathcal{O}_1(s^1, \theta^1) + \dots + \mathcal{O}_K(s^K, \theta^K)$ . In this case, the resulting  
86 block satisfies the exact same adjoint equation as that of the original learning problem (Eq. (2)) so  
87 that all gradient computation techniques herein presented may also apply *block-wise*. Alternatively,  
88 another solution to backward locking is to *estimate*  $\lambda_{k+1}$  with auxiliary modules (38).

89 **Checkpointing & reversible models.** One way to mitigate the memory problem mentioned above  
90 (at the expense of compute) is to *simultaneously* compute  $s$  and  $\lambda$ , which can be viewed as activation  
91 checkpointing (17). In models with an explicit layer hierarchy ( $\mathcal{P}_2$ ), another special instantiation  
92 of activation checkpointing is possibly when using *reversible models* (19; 26; 16; 52): a given  
93 constraint  $\mathcal{C}_{k+1}(s^{k+1}, s^k) = 0$  can be explicitly inverted as  $\mathcal{C}_{k+1}^{-1}(s^k, s^{k+1}) = 0$  such that  $s^k$  can be  
94 recomputed *backward* from layer  $s^{k+1}$  instead of being stored, or recomputed *forward* from  $s^0 = x$ .  
95 This can be achieved for instance by splitting each block state  $s^k$  into two,  $s^k = (s_a^k, s_b^k)^\top$ , and  
96 defining  $\mathcal{C}_{k+1} = (\mathcal{C}_{k+1}^a, \mathcal{C}_{k+1}^b)$ , with dedicated transformations  $f_a, f_b$ , as (26):

$$\mathcal{C}_{k+1}^a(s^{k+1}, s^k) = s_a^{k+1} - s_a^k - f_a(s_b^k, \theta_a^k), \quad \mathcal{C}_{k+1}^b(s^{k+1}, s^k) = s_b^{k+1} - s_b^k - f_b(s_a^k, \theta_b^k) \quad (7)$$

97 **Pipelining.** When dealing with a hierarchical model ( $\mathcal{P}_2$ ), the block  $\mathcal{C}_k$  may become *idle* or  
98 “forward-locked” after passing  $s^k$  to  $\mathcal{C}_{k+1}$  until next input comes in. A solution to this is to push  
99 *multiple inputs in sequence* through the blocks, allowing them to process different inputs *in parallel*,  
100 e.g.  $\mathcal{C}_k$  processes input  $x_p$  while  $\mathcal{C}_{k+1}$  processes input  $x_{p+1}$ , the same strategy applying *backwards*  
101 for the computation of the  $\lambda_k$ ’s for different inputs (37). As naive pipelining may still maintain  
102 idleness “bubbles”, more elaborate schemes have been proposed (22), for instance by allowing each  
103 block to alternate between forward pass and backward passes for different inputs (67), yet at the cost  
104 of *gradient staleness* –  $\lambda$  is computed with a different  $\theta$  than the one used to compute  $s$ . This problem  
105 can be mitigated for instances by by maintaining different weight versions at each circuit (68).

## 106 3 Forward-only learning beyond zeroth order

### 107 3.1 Energy-based (EB) models & Energy-based Learning (EBL)

108 When the constraint of the optimization problem Eq. (1) derives from an energy function  $\mathcal{K}$ , energy-  
109 based learning (EBL) refers to a family of gradient computation algorithms which implicitly estimate  
110  $s$  and  $\lambda$  using a *single* circuit (36; 32; 66; 5; 86; 78; 80). Namely, if there exists some scalar function  
111  $\mathcal{K}$  such that  $\mathcal{C} := \nabla_s \mathcal{K}$ , defining  $\mathcal{F} := \mathcal{K} + \beta \mathcal{O}$  where  $\beta \geq 0$  is some scalar value, then one can  
112 estimate  $g(\theta)$  by having the *same* circuit relax twice to equilibrium with two different values of  $\beta$   
113 and subsequent “nudged states”  $s_{\pm\beta}$  (78; 43):

$$s_{\pm\beta} : \nabla_s \mathcal{F}(\pm\beta, s_{\pm\beta}, \theta) = 0, \quad g(\theta) = \frac{1}{2\beta} (\nabla_\theta \mathcal{F}(\beta, s_\beta, \theta) - \nabla_\theta \mathcal{F}(-\beta, s_{-\beta}, \theta)) + \mathcal{O}(\beta^2) \quad (8)$$

114 Eq. (8) is in stark contrast with Eq. (2), though being mathematically equivalent through  $\lambda =$   
115  $d_\beta s_\beta|_{\beta=0}$  (78; 69): instead of estimating  $s$  and  $\lambda$  on two circuits,  $s_\beta$  and  $s_{-\beta}$  are estimated on

116 a *single* circuit through “energy minimization”, which is why  $\mathcal{C}$  is said to be *energy based* (EB).  
 117 The core intuition behind the magic of EB learning is that while error signals are usually carried  
 118 “backward” through  $\partial_s \mathcal{C}^\top$  (Eq. (2)), since we have  $\partial_s \mathcal{C}^\top = \nabla_s^2 \mathcal{K} = \partial_s \mathcal{C}$ , error signals can in this  
 119 case be equivalently carried “forward” through  $\partial_s \mathcal{C}$  through a small perturbation of  $s$  along  $\nabla_s \mathcal{O}$  of  
 120 sufficiently small  $\beta$ . Finally note that Eq. (8) is only one many variants (80) to estimate  $d_\beta(\nabla_\theta \mathcal{F})|_{\beta=0}$   
 121 (78), where there is trade-off between the number of  $s_\beta$  being evaluated and the resulting bias (94; 42).

### 122 3.2 The importance of nudging

123 **The weak & strong nudging limits.** The most current EBL setting is when  $\mathcal{C}$  describes the implicit  
 124 model itself and  $\mathcal{O}$  the cost function:  $(\mathcal{C}, \mathcal{O}) = (F = \nabla_s E, \ell)$  where we have denoted  $K = E$ . In  
 125 this case, the condition on the nudged state reads (78):

$$\nabla_s E(s_\beta, \theta) + \beta \nabla_s \ell(s_\beta, \theta) = 0 \quad (9)$$

126 We call the *nudging* factor the scalar controlling the strength of the cost function  $\ell$  in the definition of  
 127  $s_\beta$ . For this choice of  $(\mathcal{C}, \mathcal{O})$ ,  $\beta$  is the nudging factor and since theoretical guarantees of Eq. (8) hold  
 128 for  $\beta \rightarrow 0$ , it corresponds to small  $\nabla_s \ell$  perturbations, which is therefore called the *weak nudging*  
 129 *limit*. Conversely, when swapping the objective and the constraint (59; 60), i.e.  $(\mathcal{C}, \mathcal{O}) = (\nabla_s \ell, E)$ ,  
 130 which amounts to swap  $E$  and  $\ell$  inside Eq. (9), one can see that  $\beta^{-1}$  now gates the error signal  $\nabla_s \ell$ .  
 131 Therefore in this situation,  $\beta \rightarrow 0$  corresponds to *large*  $\nabla_s \ell$  perturbations, which is therefore called  
 132 the *strong nudging limit*. While it seems counter-intuitive that strong nudging solutions are relevant –  
 133 the main goal of of Eq. (1) should be to minimize the loss subject to constraints on the network energy  
 134 and not the other way around – they are global minimizers of  $\theta \rightarrow \ell(\theta, s(\theta))$  for certain choices of  $E$   
 135 (60). Strong nudging can enable learning in situations where the physical system at use is noisy by  
 136 improving the signal to noise ratio for the teaching signal (59; 42).

137 **Nudging through adiabatic oscillations.** EBL classically operates in the weak nudging regime  
 138 and estimates the error signal  $\lambda = d_\beta s_\beta|_{\beta=0}$  through a discrete, two-steps finite difference procedure,  
 139 as highlighted in Eq. (8). Another way to view the error signal is as the contour integration in the  
 140 complex plane  $\lambda = \frac{1}{T|\beta|} \int_0^T e^{i2\pi t/T} s_\beta(t)$  with  $\beta(t) := |\beta|e^{i2\pi t/T}$ . Then, provided  $T$  is sufficiently  
 141 large compared to the characteristic time constant which governs the relaxation of  $\mathcal{C}$  to equilibrium,  
 142 the error signal can be computed through slow adiabatic oscillations of the system (5; 42; 3).

143 **Nudging heuristics.** The error signal  $\beta \nabla_s \ell$  may be transmitted *instantaneously* to the rest of the  
 144 network, or with a delay. Denoting  $h(\theta, s)$  the model logits, some output controller dynamics  $u$  may  
 145 integrate the error signal inside the output layer as  $\tau_u \dot{u} + u = -\beta \nabla_o \ell(h(s, \theta))$  and feed it back to  
 146 the rest of the circuit (60). There also exists other types of nudging among the broader family of  
 147 *Contrastive Learning* (CL) algorithms. Denoting  $y$  some label, the “nudged” hidden state  $h_y$  can  
 148 be defined by hard-clamping  $y$  to the output units,  $\nabla_s E(h_y, o = y, \theta) = 0$ , while prescribing a  
 149 weight update of the same contrastive form as Eq. 8 (66). An hybrid between weak nudging and  
 150 target clamping consists of clamping the output to a *weakly nudged target*, resulting in a nudged  
 151 hidden state  $h^\beta$  defined as:  $\nabla_s E(h^\beta, o = (1 - \beta)o_* + \beta y) = 0$  where  $o_*$  corresponds to the free  
 152 state ( $\beta = 0$ ) of the output (86).

153 **Nudging synapses.** One may wonder whether neurons and parameters could play a symmetrical  
 154 role, by having them *both* evolve throughout the gradient computation phase or satisfy an equilibrium  
 155 condition. One common *post hoc* solution is to play on characteristic time constants by having  
 156 parameters slowly integrate the *instantaneous* contrastive updates prescribed by Eq. (8) (20; 87; 42;  
 157 18), which however requires the detailed knowledge of the energy function at use. Another solution  
 158 circumventing this constraint is to let *both*  $s$  and  $\theta$  equilibrate in the same energy landscape, with  
 159 control variables  $u$  acting upon  $\theta$  with strength  $\alpha$  (79). When computing  $s_{-\beta}$  (for instance),  $u$  is  
 160 adjusted to keep  $\theta$  at its current value  $\theta_t$  such that  $u_t = \theta_t + \alpha \nabla_s \mathcal{F}(-\beta, s_{-\beta}, \theta_t)$ . Therefore, keeping  
 161  $u_t$  constant when computing  $s_\beta$ ,  $\theta$  equilibrates at:  $\theta_{t+1} = u_t - \alpha \nabla_s \mathcal{F}(\beta, s_\beta, \theta_{t+1}) \approx \theta_t - \alpha \beta g(\theta)$ .

### 162 3.3 Two phases or twice as many neurons?

163 **Lagrangian reparametrization.** One may usually regard the contrastive update Eq. (8) as a  
 164 procedure performed on the *same* physical system whose state is measured at two different times with

165 two different nudging values (e.g.  $-\beta, \beta$ ), or a continuum of such. Yet, another view of Eq. (8) is to  
 166 assume that the two nudged states are computed *simultaneously by two different circuits* sharing the  
 167 same parameters. An advantage of this implementation, which has been experimentally demonstrated  
 168 with transistor-based synapses (18), is that it requires a single phase only instead of two, yet at the  
 169 expense of area and complex engineer to enable parameter sharing. A different way to approach the  
 170 same question is through the concept of *dyadic neurons* (34; 35): two variables  $s_+$  and  $s_-$  are such  
 171 that their (convex) sum encodes the model prediction and their difference the error signal if they are  
 172 critical points of the following reparametrized Lagrangian (35):

$$\tilde{\mathcal{L}}_1(s_+, s_-, \theta) := \mathcal{L}_1(\alpha s_+ + (1 - \alpha)s_-, (s_+ - s_-)/\beta), \quad (10)$$

173 where  $\mathcal{L}$  defined inside Eq. (2) and  $\alpha \in [0, 1]$ . While the resulting gradient computation algorithm  
 174 is generally a reparametrization of implicit differentiation by construction, when the constraint at  
 175 use inside Eq. (1) is energy based and  $\beta$  is sufficiently small,  $s_{\pm}$  can be simply construed as  $s_{\pm\beta}$   
 176 of Eq. (8). This can be shown by performing this reparametrization on the *relaxed* Lagrangian  $\mathcal{L}_3$   
 177 (Eq. (5)) with *fixed*  $\beta$  (35), or equivalently approximated from Eq. (10) in the limit  $\beta \rightarrow 0$ .

### 178 3.4 Applying EBL to implicit models

179 **Weak nudging.** One of the ways to cast an implicit model  $F$ , which does not explicitly derive from  
 180 an energy function, into an EB model is simply to employ the energy function  $\mathcal{K} = E = \frac{1}{2}\|F\|^2$  (60).  
 181 An important case, as one means to train feedforward nets by EP, is when  $F(s, \theta) = s - f(\theta, s)$   
 182 where  $\theta$  is typically lower block diagonal (25; 89; 62; 64). Eq. (8) for some  $\beta$ , in the weak nudging  
 183 regime ( $(\mathcal{O}, \mathcal{K}) = (\ell, E)$ ), yields in this case:

$$\begin{cases} s_\beta = f(s_\beta, \theta) + \partial_s f(s_\beta, \theta)^\top \cdot (s_\beta - f(s_\beta, \theta)) - \beta \nabla_s \ell(s_\beta, \theta), \\ g(\theta) = \frac{1}{\beta} \partial_\theta f(s_\beta, \theta)^\top \cdot (s_\beta - f(s_\beta, \theta)) + \mathcal{O}(\beta) \end{cases} \quad (11)$$

184 A nice feature of Eq. (11) is that by construction,  $s_\beta = f(s_\beta, \theta)$  when  $\beta = 0$  such that the gradient  
 185 can be estimated up to  $\mathcal{O}(\beta)$  with a *single* nudged state. However, the presence of  $\partial_s f(s_\beta, \theta)^\top$   
 186 signals the potential use analytical derivatives. To obviate this, one may instead estimate the  
 187 required derivatives through finite differences of *feedback operators* (46; 76; 57; 58; 21; 56) as:  
 188  $s_\beta \approx f(s_\beta, \theta) + g(s_\beta, \omega) - g(f(s_\beta), \omega) - \beta \nabla_s \ell(s_\beta, \theta)$  where the feedback parameters may have  
 189 to be learned so that  $\partial_s g \approx \partial_s f^\top$  (21). However a problem remaining is that both  $s_\beta$  and their  
 190 “feedforward” prediction  $f(s_\beta)$  need to be simultaneously maintained to implement Eq. (11). A  
 191 solution is to use auxiliary variables  $\psi$  which track the feedforward activity such that  $\psi_\beta \approx f_\psi(s_\beta, \theta)$   
 192 and  $s_\beta \approx f(s_\beta, \theta) + g(s_\beta, \omega) - g(\psi_\beta, \omega) - \beta \nabla_s \ell(s_\beta, \theta)$  (76).

193 **Recovering a least control problem in the strong nudging regime.** When in the strong nudging  
 194 regime, the optimization problem Eq. (1) using  $E = \frac{1}{2}\|F\|^2 = \mathcal{O}$  along with  $\ell = \mathcal{K}$  acquires a  
 195 very intriguing meaning, though the loss is not picked as the objective. Defining  $\psi_\beta := -F(s_\beta, \theta)$ ,  
 196 the learning problem amounts to find the set of parameters which minimizes the amount of *control*  
 197  $\mathcal{O} = \frac{1}{2}\|\psi_\beta\|^2$  such that the controlled equilibrium  $\psi_\beta + F(s_\beta, \theta) = 0$  minimizes the loss  $\ell$ , which  
 198 works in practice on feedforward and implicit models more broadly (60). However, this interpretation  
 199 is tied to the choice  $E = \frac{1}{2}\|F\|^2$  and further investigation is needed to extend the applicability of  
 200 strong nudging and its connection to control theory beyond this choice of energy function.

### 201 3.5 Handling nonlinearity inside EB models

202 When not handled cautiously, the application of Eq. 8 to implicit models yields activation derivatives  
 203 (Eq. (11)). While neuroscientists, seeking for biologically plausible learning theories, gave birth  
 204 to the aforementioned approximations of Eq. (11) and many others, these remain impractical for  
 205 hardware-friendly learning. One *ad hoc* solution used in practice, when the inverse of the activation  
 206 function is continuously invertible, is to separate the linear and nonlinear contribution inside the  
 207 energy function as  $E = G + U$  where  $G$  is defined such that  $\nabla_s G = \sigma^{-1}$  (36; 92; 34; 80). A  
 208 more principled solution to handle  $\sigma$  is to instead treat it as a *constraint*  $\mathcal{C}_\sigma$  on the set of feasible  
 209 configurations over which energy minimization is performed, which descent steps are *hard* projected  
 210 onto(77). A “relaxed” version of this approach, somewhat bridging with the former solution, is  
 211 through the lens of *mirror descent* using the Bregman divergence associated with  $G$  (2):  $s_{t+1} =$   
 212  $\arg \min\{s^\top \cdot \nabla_s E(s) + \frac{1}{\eta} \Delta_G(s, s_t)\}$  with  $\Delta_G(x, y) := G(x) - G(y) + \nabla G(y)^\top \cdot (x - y)$ . Note  
 213 that many other fixed point algorithms could be use for the same purpose (13).

### 214 3.6 Casting feedforward nets into hierarchical energy-based models

215 **An important example.** Another way to cast feedforward nets training by EBL algorithms is to  
 216 break them into multiple, hierarchical energy-based constraints based on Eq. (3) (92). The following  
 217 example, which shows how to do so, is also a good illustration of why the problems  $\mathcal{P}_1$  (Eq. (1))  
 218 and  $\mathcal{P}_2$  (Eq. (3)) are generally not equivalent. Taking  $E := \sum_{k=1}^K E_k$  with  $E_k(s^k, s^{k-1}, \theta^k) :=$   
 219  $G(s^k) - s^{k\top} \cdot \theta^k \cdot s^{k-1}$ , we consider  $\mathcal{P}_1$  with  $\mathcal{C} = \nabla_s E$  and  $\mathcal{P}_2$  with  $\mathcal{C}_k = \nabla_{s^k} E^k$ . In this case,  $s^k$   
 220 is a *single layer* and not a block of several layers. Setting the constraint to zero in these two cases  
 221 yields:

$$\mathcal{P}_1 : s^k = \sigma \left( \theta^k \cdot s^{k-1} + \theta^{k+1\top} \cdot s^{k+1} \right), \quad \mathcal{P}_2 : s^k = \sigma \left( \theta^k \cdot s^{k-1} \right) \quad (12)$$

222 In one case, the model retained at inference time is a *Hopfield model* and in the other one a *feedforward*  
 223 model. This difference is instrumental to understand why works revolving around “contrastive  
 224 learning” may actually apply to feedforward models through the  $\mathcal{P}_2$  problem formulation (34).  
 225 Likewise, the nudged state defined by Eq. (8) reads for any layer *until penultimate* the same as  
 226 Eq. (12) for problem  $\mathcal{P}_1$  and as follows for  $\mathcal{P}_2$ , taking  $\beta$  sufficiently small (92; 69):

$$\mathcal{P}_2 : s_{\pm\beta}^k = \sigma \left( \theta^k \cdot s^{k-1} \pm \beta \theta^{k+1\top} \cdot d_{\beta} s_{\beta}^{k+1} \Big|_{\beta=0} \right) \approx \sigma \left( \theta^k \cdot s^{k-1} \pm \theta^{k+1\top} \cdot \left( s_{\beta}^{k+1} - s_{-\beta}^{k+1} \right) / 2 \right), \quad (13)$$

227 Note that many other choices of energy functions are possible in order to learn feedforward nets as a  
 228 hierarchical energy-based model (15; 47; 29; 92)

229 **Breaking the forward–backward synchrony.** As seen from Eq. (13), computing  $s_{\pm\beta}^k$  still requires  
 230 the storage of  $s^{k-1}$ . A trick is to notice that for sufficiently small  $\beta$ ,  $s^k \approx (s_{\beta}^k + s_{-\beta}^k) / 2 := \bar{s}_{\beta}$  such  
 231 that  $s_{\beta}^k \approx \sigma \left( \theta^k \cdot \bar{s}_{\beta} \pm \theta^{k+1\top} \cdot \Delta s_{\beta}^{k+1} \right)$  with  $\Delta s_{\beta}^{k+1} := s_{\beta}^{k+1} - s_{-\beta}^{k+1}$ . Implemented this way, the  
 232 gradient corresponding to problem  $\mathcal{P}^2$  with the above choice of energy function can be computed by  
 233 *simultaneously* solving for  $s_{\beta}^k$ ’s and  $s_{-\beta}^k$ ’s (34). This implementation can also be derived by picking  
 234 hierarchical constraints (Eq. (3)), relaxing these (Eq. (5)) and reparametrizing them (Eq. (10)) (35).

## 235 4 Looking ahead & take-aways

236 **Forward-only learning beyond first order.** One may wonder if (and how) *second-order opti-*  
 237 *mization* could be emulated using a *single* circuit. A second-order update generally prescribes  
 238  $\Delta \theta_{2\text{nd}} \propto \mathbb{E}_{x \sim \mathcal{D}} [C(x, \theta)]^{-1} \cdot \mathbb{E}_{x \sim \mathcal{D}} [g(x, \theta)]$  where  $\mathcal{C}$  is a curvature matrix, i.e. the loss *Hessian*  
 239 or p.s.d. approximations thereof (53). As such, computing  $\theta_{2\text{nd}}$  generally subsumes (*Hessian-free*  
 240 approaches being a notable exception (55)): i) computing  $g$ , ii) computing  $C$ , iii) inverting  $C$ , iv)  
 241 computing  $\Delta \theta_{2\text{nd}}$ . However, an interesting insight arises when  $|\mathcal{D}| = 1$  (unit batch size), picking  
 242 the Fisher information curvature matrix  $C = (\partial_{\theta} s^{\top} \cdot \partial_{\theta} s)$  (54) and assuming  $\mathcal{O}$  does not explicitly  
 243 depend on  $\theta$  for simplicity. Then, denoting  $A^{\dagger}$  the *Moore-Penrose pseudo-inverse* of a matrix  $A$ :

$$\Delta \theta_{2\text{nd}} \propto (\partial_{\theta} s^{\top} \cdot \partial_{\theta} s)^{-1} \cdot \partial_{\theta} s^{\top} \cdot \nabla_s \mathcal{O} = \partial_{\theta} s^{\dagger} \cdot \nabla_s \mathcal{O}. \quad (14)$$

244 Therefore, moving to second order optimization here amounts to route the error signal  $\nabla_s \mathcal{O}$  backward  
 245 through the system’s *inverse* rather than its adjoint. While this viewpoint has motivated several  
 246 implementations of *self-invertible* (rather than *self-adjoint*, or equivalently energy-based) circuits (46;  
 247 57; 11; 58), it may not straightforwardly generalize to the mini-batch regime ( $|\mathcal{D}| > 1$ ).

248 **Temporal data?** One way to leverage equilibrium techniques towards learning from sequences is  
 249 to treat the neuron *velocity*  $\dot{s}$  as a variable of its own inside the energy function (30; 82).

250 **Conclusion.** This “cookbook” is an attempt to complement recent algorithm unification papers  
 251 (92; 93; 63; 72) using pragmatic, ingredient-based logics. For instance, “Predictive coding” (89)  
 252 amounts to apply Eq. (8) with a specific choice of energy function (3.4) and nudging scheme (3.2).  
 253 Likewise, the “Forward-forward” algorithm (31) is a greedy learning strategy, with each block  
 254 comprising a single layer, using local self-supervised loss (2.3). Conversely, one may want to split  
 255 an architecture into blocks (Eq. (3)), apply standard implicit differentiation in some (2.2, §1), ZO  
 256 in others (2.2, §1), along with a pipelining mechanism (2.3) depending on the constraints at hand,  
 257 which does not correspond to any known algorithm in the literature.

## References

- 258
- 259 [1] M. Akrouf, C. Wilson, P. Humphreys, T. Lillicrap, and D. B. Tweed. Deep learning without  
260 weight transport. *Advances in neural information processing systems*, 32, 2019.
- 261 [2] E. Amid, R. Anil, and M. Warmuth. Locoprop: Enhancing backprop via local loss optimization.  
262 In *International Conference on Artificial Intelligence and Statistics*, pages 9626–9642. PMLR,  
263 2022.
- 264 [3] V. R. Anisetti, A. Kandala, B. Scellier, and J. Schwarz. Frequency propagation: Multimechanism  
265 learning in nonlinear physical networks. *Neural Computation*, 36(4):596–620, 2024.
- 266 [4] S. Bai, J. Z. Kolter, and V. Koltun. Deep equilibrium models. *Advances in neural information  
267 processing systems*, 32, 2019.
- 268 [5] P. Baldi and F. Pineda. Contrastive learning and neural oscillations. *Neural computation*, 3(4):  
269 526–545, 1991.
- 270 [6] S. Bartunov, A. Santoro, B. Richards, L. Marris, G. E. Hinton, and T. Lillicrap. Assessing the  
271 scalability of biologically-motivated deep learning algorithms and architectures. *Advances in  
272 neural information processing systems*, 31, 2018.
- 273 [7] A. G. Baydin, B. A. Pearlmutter, D. Syme, F. Wood, and P. Torr. Gradients without backpropa-  
274 gation. *arXiv preprint arXiv:2202.08587*, 2022.
- 275 [8] E. Belilovsky, M. Eickenberg, and E. Oyallon. Greedy layerwise learning can scale to imagenet.  
276 In *International conference on machine learning*, pages 583–593. PMLR, 2019.
- 277 [9] E. Belilovsky, L. Leconte, L. Caccia, M. Eickenberg, and E. Oyallon. Decoupled greedy  
278 learning of cnns for synchronous and asynchronous distributed learning. *arXiv preprint  
279 arXiv:2106.06401*, 2021.
- 280 [10] B. M. Bell and J. V. Burke. Algorithmic differentiation of implicit functions and optimal values.  
281 In *Advances in automatic differentiation*, pages 67–77. Springer, 2008.
- 282 [11] Y. Bengio. Deriving differential target propagation from iterating approximate inverses. *arXiv  
283 preprint arXiv:2007.15139*, 2020.
- 284 [12] D. P. Bertsekas. *Constrained optimization and Lagrange multiplier methods*. Academic press,  
285 2014.
- 286 [13] M. Blondel, Q. Berthet, M. Cuturi, R. Frostig, S. Hoyer, F. Llinares-López, F. Pedregosa,  
287 and J.-P. Vert. Efficient and modular implicit differentiation. *Advances in neural information  
288 processing systems*, 35:5230–5242, 2022.
- 289 [14] G. Bouvier, J. Aljadeff, C. Clopath, C. Bimbard, J. Ranft, A. Blot, J.-P. Nadal, N. Brunel,  
290 V. Hakim, and B. Barbour. Cerebellar learning using perturbations. *Elife*, 7:e31599, 2018.
- 291 [15] M. Carreira-Perpinan and W. Wang. Distributed optimization of deeply nested systems. In  
292 *Artificial Intelligence and Statistics*, pages 10–19. PMLR, 2014.
- 293 [16] B. Chang, L. Meng, E. Haber, L. Ruthotto, D. Begert, and E. Holtham. Reversible architectures  
294 for arbitrarily deep residual neural networks. In *Proceedings of the AAAI conference on artificial  
295 intelligence*, volume 32, 2018.
- 296 [17] T. Chen, B. Xu, C. Zhang, and C. Guestrin. Training deep nets with sublinear memory cost.  
297 *arXiv preprint arXiv:1604.06174*, 2016.
- 298 [18] S. Dillavou, B. D. Beyer, M. Stern, M. Z. Miskin, A. J. Liu, and D. J. Durian. Machine learning  
299 without a processor: Emergent learning in a nonlinear electronic metamaterial. *arXiv preprint  
300 arXiv:2311.00537*, 2023.
- 301 [19] L. Dinh, D. Krueger, and Y. Bengio. Nice: Non-linear independent components estimation.  
302 *arXiv preprint arXiv:1410.8516*, 2014.

- 303 [20] M. Ernoult, J. Grollier, D. Querlioz, Y. Bengio, and B. Scellier. Equilibrium propagation with  
304 continual weight updates. *arXiv preprint arXiv:2005.04168*, 2020.
- 305 [21] M. M. Ernoult, F. Normandin, A. Moudgil, S. Spinney, E. Belilovsky, I. Rish, B. Richards,  
306 and Y. Bengio. Towards scaling difference target propagation by learning backprop targets. In  
307 *International Conference on Machine Learning*, pages 5968–5987. PMLR, 2022.
- 308 [22] S. Fan, Y. Rong, C. Meng, Z. Cao, S. Wang, Z. Zheng, C. Wu, G. Long, J. Yang, L. Xia, et al.  
309 Dapple: A pipelined data parallel approach for training large models. In *Proceedings of the*  
310 *26th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pages  
311 431–445, 2021.
- 312 [23] I. R. Fiete, M. S. Fee, and H. S. Seung. Model of birdsong learning based on gradient estimation  
313 by dynamic perturbation of neural conductances. *Journal of neurophysiology*, 98(4):2038–2057,  
314 2007.
- 315 [24] L. Fournier, S. Rivaud, E. Belilovsky, M. Eickenberg, and E. Oyallon. Can forward gradient  
316 match backpropagation? In *International Conference on Machine Learning*, pages 10249–  
317 10264. PMLR, 2023.
- 318 [25] K. Friston and S. Kiebel. Predictive coding under the free-energy principle. *Philosophical*  
319 *transactions of the Royal Society B: Biological sciences*, 364(1521):1211–1221, 2009.
- 320 [26] A. N. Gomez, M. Ren, R. Urtasun, and R. B. Grosse. The reversible residual network: Back-  
321 propagation without storing activations. *Advances in neural information processing systems*, 30,  
322 2017.
- 323 [27] A. N. Gomez, O. Key, K. Perlin, S. Gou, N. Frosst, J. Dean, and Y. Gal. Interlocking backprop-  
324 agation: Improving depthwise model-parallelism. *Journal of Machine Learning Research*, 23  
325 (171):1–28, 2022.
- 326 [28] A. Griewank and A. Walther. *Evaluating derivatives: principles and techniques of algorithmic*  
327 *differentiation*. SIAM, 2008.
- 328 [29] F. Gu, A. Askari, and L. El Ghaoui. Fenchel lifted networks: A lagrange relaxation of neural  
329 network training. In *International Conference on Artificial Intelligence and Statistics*, pages  
330 3362–3371. PMLR, 2020.
- 331 [30] P. Haider, B. Ellenberger, L. Kriener, J. Jordan, W. Senn, and M. A. Petrovici. Latent equilibrium:  
332 A unified learning theory for arbitrarily fast computation with arbitrarily slow neurons. *Advances*  
333 *in neural information processing systems*, 34:17839–17851, 2021.
- 334 [31] G. Hinton. The forward-forward algorithm: Some preliminary investigations. *arXiv preprint*  
335 *arXiv:2212.13345*, 2022.
- 336 [32] G. E. Hinton, T. J. Sejnowski, and D. H. Ackley. *Boltzmann machines: Constraint satisfaction*  
337 *networks that learn*. Carnegie-Mellon University, Department of Computer Science Pittsburgh,  
338 PA, 1984.
- 339 [33] N. Hiratani, Y. Mehta, T. Lillicrap, and P. E. Latham. On the stability and scalability of node  
340 perturbation learning. *Advances in Neural Information Processing Systems*, 35:31929–31941,  
341 2022.
- 342 [34] R. Høier, D. Staudt, and C. Zach. Dual propagation: Accelerating contrastive hebbian learning  
343 with dyadic neurons. In *International Conference on Machine Learning*, pages 13141–13156.  
344 PMLR, 2023.
- 345 [35] R. K. Høier and C. Zach. Two tales of single-phase contrastive hebbian learning. *arXiv preprint*  
346 *arXiv:2402.08573*, 2024.
- 347 [36] J. J. Hopfield. Neurons with graded response have collective computational properties like those  
348 of two-state neurons. *Proceedings of the national academy of sciences*, 81(10):3088–3092,  
349 1984.



- 350 [37] Y. Huang, Y. Cheng, A. Bapna, O. Firat, D. Chen, M. Chen, H. Lee, J. Ngiam, Q. V. Le, Y. Wu,  
351 et al. Gpipe: Efficient training of giant neural networks using pipeline parallelism. *Advances in*  
352 *neural information processing systems*, 32, 2019.
- 353 [38] M. Jaderberg, W. M. Czarnecki, S. Osindero, O. Vinyals, A. Graves, D. Silver, and  
354 K. Kavukcuoglu. Decoupled neural interfaces using synthetic gradients. In *International*  
355 *conference on machine learning*, pages 1627–1635. PMLR, 2017.
- 356 [39] J. Kendall, R. Pantone, K. Manickavasagam, Y. Bengio, and B. Scellier. Training end-to-end  
357 analog neural networks with equilibrium propagation. *arXiv preprint arXiv:2006.01981*, 2020.
- 358 [40] J. Kornfeld, M. Januszewski, P. Schubert, V. Jain, W. Denk, and M. S. Fee. An anatomical  
359 substrate of credit assignment in reinforcement learning. *BioRxiv*, pages 2020–02, 2020.
- 360 [41] S. G. Krantz and H. R. Parks. *The implicit function theorem: history, theory, and applications*.  
361 Springer Science & Business Media, 2002.
- 362 [42] A. Laborieux and F. Zenke. Holomorphic equilibrium propagation computes exact gradients  
363 through finite size oscillations. *Advances in neural information processing systems*, 35:12950–  
364 12963, 2022.
- 365 [43] A. Laborieux, M. Ernoult, B. Scellier, Y. Bengio, J. Grollier, and D. Querlioz. Scaling equilib-  
366 rium propagation to deep convnets by drastically reducing its gradient estimator bias. *Frontiers*  
367 *in neuroscience*, 15:633674, 2021.
- 368 [44] B. J. Lansdell, P. R. Prakash, and K. P. Kording. Learning to solve the credit assignment  
369 problem. *arXiv preprint arXiv:1906.00889*, 2019.
- 370 [45] J. Launay, I. Poli, F. Boniface, and F. Krzakala. Direct feedback alignment scales to modern  
371 deep learning tasks and architectures. *Advances in neural information processing systems*, 33:  
372 9346–9360, 2020.
- 373 [46] D.-H. Lee, S. Zhang, A. Fischer, and Y. Bengio. Difference target propagation. In *Machine*  
374 *Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD 2015,*  
375 *Porto, Portugal, September 7-11, 2015, Proceedings, Part I 15*, pages 498–515. Springer, 2015.
- 376 [47] J. Li, C. Fang, and Z. Lin. Lifted proximal operator machines. In *Proceedings of the AAAI*  
377 *Conference on Artificial Intelligence*, volume 33, pages 4181–4188, 2019.
- 378 [48] T. P. Lillicrap, D. Cownden, D. B. Tweed, and C. J. Akerman. Random feedback weights  
379 support learning in deep neural networks. *arXiv preprint arXiv:1411.0247*, 2014.
- 380 [49] S. Liu, P.-Y. Chen, B. Kailkhura, G. Zhang, A. O. Hero III, and P. K. Varshney. A primer  
381 on zeroth-order optimization in signal processing and machine learning: Principals, recent  
382 advances, and applications. *IEEE Signal Processing Magazine*, 37(5):43–54, 2020.
- 383 [50] S. Löwe, P. O’Connor, and B. Veeling. Putting an end to end-to-end: Gradient-isolated learning  
384 of representations. *Advances in neural information processing systems*, 32, 2019.
- 385 [51] S. Malladi, T. Gao, E. Nichani, A. Damian, J. D. Lee, D. Chen, and S. Arora. Fine-tuning  
386 language models with just forward passes. *Advances in Neural Information Processing Systems*,  
387 36:53038–53075, 2023.
- 388 [52] K. Mangalam, H. Fan, Y. Li, C.-Y. Wu, B. Xiong, C. Feichtenhofer, and J. Malik. Reversible  
389 vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and*  
390 *Pattern Recognition*, pages 10830–10840, 2022.
- 391 [53] J. Martens. New insights and perspectives on the natural gradient method. *Journal of Machine*  
392 *Learning Research*, 21(146):1–76, 2020.
- 393 [54] J. Martens and R. Grosse. Optimizing neural networks with kronecker-factored approximate  
394 curvature. In *International conference on machine learning*, pages 2408–2417. PMLR, 2015.

- 395 [55] J. Martens and I. Sutskever. Learning recurrent neural networks with hessian-free optimization.  
396 In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pages  
397 1033–1040, 2011.
- 398 [56] K. Max, L. Kriener, G. Pineda García, T. Nowotny, I. Jaras, W. Senn, and M. A. Petrovici.  
399 Learning efficient backprojections across cortical hierarchies in real time. *Nature Machine*  
400 *Intelligence*, pages 1–12, 2024.
- 401 [57] A. Meulemans, F. Carzaniga, J. Suykens, J. Sacramento, and B. F. Grewe. A theoretical  
402 framework for target propagation. *Advances in Neural Information Processing Systems*, 33:  
403 20024–20036, 2020.
- 404 [58] A. Meulemans, M. Tristany Farinha, J. García Ordóñez, P. Vilimelis Aceituno, J. Sacramento,  
405 and B. F. Grewe. Credit assignment in neural networks through deep feedback control. *Advances*  
406 *in Neural Information Processing Systems*, 34:4674–4687, 2021.
- 407 [59] A. Meulemans, M. T. Farinha, M. R. Cervera, J. Sacramento, and B. F. Grewe. Minimizing  
408 control for credit assignment with strong feedback. In *International Conference on Machine*  
409 *Learning*, pages 15458–15483. PMLR, 2022.
- 410 [60] A. Meulemans, N. Zucchet, S. Kobayashi, J. Von Oswald, and J. Sacramento. The least-control  
411 principle for local learning at equilibrium. *Advances in Neural Information Processing Systems*,  
412 35:33603–33617, 2022.
- 413 [61] T. Miconi. Biologically plausible learning in recurrent neural networks reproduces neural  
414 dynamics observed during cognitive tasks. *Elife*, 6:e20899, 2017.
- 415 [62] B. Millidge, A. Seth, and C. L. Buckley. Predictive coding: a theoretical and experimental  
416 review. *arXiv preprint arXiv:2107.12979*, 2021.
- 417 [63] B. Millidge, Y. Song, T. Salvatori, T. Lukasiewicz, and R. Bogacz. Backpropagation at the  
418 infinitesimal inference limit of energy-based models: Unifying predictive coding, equilibrium  
419 propagation, and contrastive hebbian learning. *arXiv preprint arXiv:2206.02629*, 2022.
- 420 [64] B. Millidge, A. Tschantz, and C. L. Buckley. Predictive coding approximates backprop along  
421 arbitrary computation graphs. *Neural Computation*, 34(6):1329–1368, 2022.
- 422 [65] T. H. Moskovitz, A. Litwin-Kumar, and L. Abbott. Feedback alignment in deep convolutional  
423 networks. *arXiv preprint arXiv:1812.06488*, 2018.
- 424 [66] J. R. Movellan. Contrastive hebbian learning in the continuous hopfield model. In *Connectionist*  
425 *models*, pages 10–17. Elsevier, 1991.
- 426 [67] D. Narayanan, A. Harlap, A. Phanishayee, V. Seshadri, N. R. Devanur, G. R. Ganger, P. B.  
427 Gibbons, and M. Zaharia. Pipedream: Generalized pipeline parallelism for dnn training. In  
428 *Proceedings of the 27th ACM symposium on operating systems principles*, pages 1–15, 2019.
- 429 [68] D. Narayanan, A. Phanishayee, K. Shi, X. Chen, and M. Zaharia. Memory-efficient pipeline-  
430 parallel dnn training. In *International Conference on Machine Learning*, pages 7937–7947.  
431 PMLR, 2021.
- 432 [69] T. Nest and M. Ernout. Towards training digitally-tied analog blocks via hybrid gradient  
433 computation. *arXiv preprint arXiv:2409.03306*, 2024.
- 434 [70] A. Nøkland. Direct feedback alignment provides learning in deep neural networks. *Advances in*  
435 *neural information processing systems*, 29, 2016.
- 436 [71] A. Nøkland and L. H. Eidnes. Training neural networks with local error signals. In *International*  
437 *conference on machine learning*, pages 4839–4850. PMLR, 2019.
- 438 [72] A. Ororbiala, A. Mali, A. Kohan, B. Millidge, and T. Salvatori. A review of neuroscience-inspired  
439 machine learning. *arXiv preprint arXiv:2403.18929*, 2024.
- 440 [73] J. V. Outrata. A note on the usage of nondifferentiable exact penalties in some special optimiza-  
441 tion problems. *Kybernetika*, 24(4):251–258, 1988.

- 442 [74] M. Ren, S. Kornblith, R. Liao, and G. Hinton. Scaling forward gradient with local losses. *arXiv*  
443 *preprint arXiv:2210.03310*, 2022.
- 444 [75] Y. Ruan, Y. Xiong, S. Reddi, S. Kumar, and C.-J. Hsieh. Learning to learn by zeroth-order  
445 oracle. *arXiv preprint arXiv:1910.09464*, 2019.
- 446 [76] J. Sacramento, R. Ponte Costa, Y. Bengio, and W. Senn. Dendritic cortical microcircuits  
447 approximate the backpropagation algorithm. *Advances in neural information processing systems*,  
448 31, 2018.
- 449 [77] B. Scellier. A fast algorithm to simulate nonlinear resistive networks. *arXiv preprint*  
450 *arXiv:2402.11674*, 2024.
- 451 [78] B. Scellier and Y. Bengio. Equilibrium propagation: Bridging the gap between energy-based  
452 models and backpropagation. *Frontiers in computational neuroscience*, 11:24, 2017.
- 453 [79] B. Scellier, S. Mishra, Y. Bengio, and Y. Ollivier. Agnostic physics-driven deep learning. *arXiv*  
454 *preprint arXiv:2205.15021*, 2022.
- 455 [80] B. Scellier, M. Ernoult, J. Kendall, and S. Kumar. Energy-based learning algorithms for analog  
456 computing: a comparative study. *Advances in Neural Information Processing Systems*, 36,  
457 2024.
- 458 [81] G. Scheithauer. Jorge nocedal and stephen j. wright: Numerical optimization, springer series in  
459 operations research, 1999, isbn 0-387-98793-2 in the preface the authors state... our goal in this  
460 book is to give a comprehensive description of the most powerful, state-of-the-art, techniques  
461 for solving continuous optimization problems. by presenting the motivating ideas for each.
- 462 [82] W. Senn, D. Dold, A. F. Kungl, B. Ellenberger, J. Jordan, Y. Bengio, J. Sacramento, and M. A.  
463 Petrovici. A neuronal least-action principle for real-time learning in cortical circuits. *BioRxiv*,  
464 pages 2023–03, 2023.
- 465 [83] S. A. Siddiqui, D. Krueger, Y. LeCun, and S. Deny. Blockwise self-supervised learning at scale.  
466 *arXiv preprint arXiv:2302.01647*, 2023.
- 467 [84] D. Silver, A. Goyal, I. Danihelka, M. Hessel, and H. van Hasselt. Learning by directional  
468 gradient descent. In *International Conference on Learning Representations*, 2021.
- 469 [85] J. C. Spall. Multivariate stochastic approximation using a simultaneous perturbation gradient  
470 approximation. *IEEE transactions on automatic control*, 37(3):332–341, 1992.
- 471 [86] M. Stern, D. Hexner, J. W. Rocks, and A. J. Liu. Supervised learning in physical networks:  
472 From machine learning to learning machines. *Physical Review X*, 11(2):021045, 2021.
- 473 [87] M. Stern, S. Dillavou, M. Z. Miskin, D. J. Durian, and A. J. Liu. Physical learning beyond the  
474 quasistatic limit. *Physical Review Research*, 4(2):L022037, 2022.
- 475 [88] Y. Wang, Z. Ni, S. Song, L. Yang, and G. Huang. Revisiting locally supervised learning: an  
476 alternative to end-to-end training. *arXiv preprint arXiv:2101.10832*, 2021.
- 477 [89] J. C. Whittington and R. Bogacz. An approximation of the error backpropagation algorithm in a  
478 predictive coding network with local hebbian synaptic plasticity. *Neural computation*, 29(5):  
479 1229–1262, 2017.
- 480 [90] W. Xiao, H. Chen, Q. Liao, and T. Poggio. Biologically-plausible learning algorithms can scale  
481 to large datasets. *arXiv preprint arXiv:1811.03567*, 2018.
- 482 [91] Y. Xiong, M. Ren, and R. Urtasun. Loco: Local contrastive representation learning. *Advances*  
483 *in neural information processing systems*, 33:11142–11153, 2020.
- 484 [92] C. Zach. Bilevel programs meet deep learning: A unifying view on inference learning methods.  
485 *arXiv preprint arXiv:2105.07231*, 2021.
- 486 [93] N. Zucchet and J. Sacramento. Beyond backpropagation: bilevel optimization through implicit  
487 differentiation and equilibrium propagation. *Neural Computation*, 34(12):2309–2346, 2022.
- 488 [94] N. Zucchet, S. Schug, J. Von Oswald, D. Zhao, and J. Sacramento. A contrastive rule for  
489 meta-learning. *Advances in neural information processing systems*, 35:25921–25936, 2022.