

## 565 A Theoretical proofs

566 In this appendix we detail the proofs of the theoretical results in the body text.

### 567 A.1 Complex analysis background

568 We recall here the minimal complex analysis background required to appreciate the theoretical results  
 569 of this work. In the following, we recall the definitions of holomorphic and Wirtinger derivatives, the  
 570 Cauchy-Riemann equations and the Cauchy formulas. We refer the reader to Chapter 4 of [S1] for  
 571 proofs as well as an excellent introduction to complex analysis.

572 **Definition 1** (Holomorphic function). *Let  $U$  be an open set of  $\mathbb{C}$  and  $f : z \in U \mapsto f(z) \in \mathbb{C}$  a*  
 573 *function.  $f$  is holomorphic at  $a \in U$  if the limit*

$$\lim_{z \rightarrow a} \frac{f(z) - f(a)}{z - a}$$

574 *exists. This limit is then noted  $f'(a)$ .  $f$  is holomorphic on  $U$  if it is holomorphic  $\forall a \in U$ .*

575 Though this definition looks like the definition of differentiability in  $\mathbb{R}$ , it brings constraints on the  
 576 underlying function  $\tilde{f} : (x, y) \in \mathbb{R}^2 \mapsto (\text{Re}(f(x + iy)), \text{Im}(f(x + iy)))$ . The added constraints are  
 577 the Cauchy-Riemann equations, which can be compactly written after defining Wirtinger derivatives:

578 **Definition 2** (Wirtinger derivatives). *Noting  $\partial/\partial x$  and  $\partial/\partial y$  the usual partial derivatives in  $\mathbb{R}^2$ ,*  
 579 *the Wirtinger derivatives are defined by:*

$$\frac{\partial}{\partial z} := \frac{1}{2} \left( \frac{\partial}{\partial x} - i \frac{\partial}{\partial y} \right), \quad \frac{\partial}{\partial \bar{z}} := \frac{1}{2} \left( \frac{\partial}{\partial x} + i \frac{\partial}{\partial y} \right).$$

580 In this way,  $z$  and its complex conjugate  $\bar{z}$  can be thought of as independent variables. We can then  
 581 state the Cauchy-Riemann equations as:

582 **Theorem 2** (Cauchy-Riemann equations). *if  $f$  is holomorphic at  $a \in U$ , then:*

$$\frac{\partial f}{\partial z}(a) = f'(a), \quad \frac{\partial f}{\partial \bar{z}}(a) = 0. \quad (11)$$

583 These constraints ensure that  $f$  is locally expandable everywhere in  $U$  into a converging power series.  
 584 In particular, it is differentiable at any order and the derivatives can be computed with the:

585 **Theorem 3** (Cauchy formulas). *Let  $f$  be holomorphic on  $U$ , let  $\gamma$  be any piece-wise continuously*  
 586 *differentiable closed curve in  $U$  going around  $a \in U$  once and counterclockwise, then:*

$$f^{(n)}(a) = \frac{n!}{2i\pi} \oint_{\gamma} \frac{f(z)}{(z - a)^{n+1}} dz. \quad (12)$$

### 587 A.2 Proof of Lemma 1

588 Here, we give a more detailed proof of holomorphic EP. We recall the:

589 **Lemma 1** (Holomorphic Equilibrium Propagation). *Let  $F$  be a scalar function governing the*  
 590 *dynamics, so that the holomorphic implicit function theorem can be applied to the fixed point*  
 591 *equation  $\partial_s F(\theta, s_0^*, 0) = 0$ , then the gradient formula of equilibrium propagation (Eq. (2)) holds in*  
 592 *the sense of complex differentiation.*

593 *Proof.* We first detail precisely the set of equations on which the holomorphic implicit function  
 594 theorem is applied. At the free fixed point  $(\theta = \theta_0, \beta = 0)$  that which exists by assumption, we have  
 595 the following set of equations:

$$\frac{\partial F}{\partial s_j}(\theta_0, s_0^*, 0) = 0, \quad 1 \leq j \leq n,$$

where  $n$  is the number of units in the system. The functions  $\partial_{s_j} F$  are holomorphic by assumption. If we further assume that the Hessian of  $F$  with respect to  $\mathbf{s}$  is invertible in  $(\boldsymbol{\theta}_0, \mathbf{s}_0^*, 0)$ , i.e.:

$$\det \left( \frac{\partial^2 F}{\partial s_i \partial s_j}(\boldsymbol{\theta}_0, \mathbf{s}_0^*, 0) \right)_{i,j} \neq 0,$$

then the holomorphic version of the implicit function theorem [S2] can be applied and there exists an open neighbourhood of  $(\boldsymbol{\theta} = \boldsymbol{\theta}_0, \beta = 0)$  in the complex domain where the implicit map  $(\boldsymbol{\theta}, \beta) \mapsto \mathbf{s}_{\boldsymbol{\theta}, \beta}^*$  is holomorphic, and where the fixed point equations hold:

$$\frac{\partial F}{\partial \mathbf{s}}(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) = 0.$$

At such fixed points, we have that the total derivatives of  $F$  with respect to either  $\beta$  or  $\boldsymbol{\theta}$  are equal to the partial derivatives, which can be seen by applying the chain rule of complex differentiation using Wirtinger derivatives. There are now in principle three contributions to the total derivative of  $F$  with respect to  $\beta$ :

$$\frac{dF}{d\beta}(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) = \frac{\partial F}{\partial \beta}(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) + \underbrace{\frac{\partial F}{\partial \mathbf{s}} \cdot \frac{\partial \mathbf{s}}{\partial \beta}(\boldsymbol{\theta}, \beta)}_{=0 \text{ at a fixed point}} + \underbrace{\frac{\partial F}{\partial \bar{\mathbf{s}}} \cdot \frac{\partial \bar{\mathbf{s}}}{\partial \beta}(\boldsymbol{\theta}, \beta)}_{=0 \text{ by Cauchy-Riemann (Eq. (11))}}, \quad (13)$$

where  $\bar{\mathbf{s}}$  denotes the complex conjugate of  $\mathbf{s}$ . At the fixed point however, the second term on the right hand side cancels by definition. The third term is zero because  $F$  is holomorphic, i.e., its derivative with respect to the conjugate variable  $\bar{\mathbf{s}}$  is zero according to the Cauchy-Riemann condition [S1]. The same argument holds for the total derivative with respect to  $\boldsymbol{\theta}$ .

Finally, the cross-derivatives of  $F$  with respect to complex  $\beta$  and  $\boldsymbol{\theta}$  can be exchanged, which is a consequence of the Schwarz theorem applied to the function  $(\boldsymbol{\theta}, \beta) \mapsto G(\boldsymbol{\theta}, \beta) := F(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta)$ . Therefore we have that:

$$\begin{aligned} \frac{\partial^2 G}{\partial \beta \partial \boldsymbol{\theta}}(\boldsymbol{\theta}, \beta) &= \frac{\partial^2 G}{\partial \boldsymbol{\theta} \partial \beta}(\boldsymbol{\theta}, \beta), \\ \frac{d}{d\beta} \frac{d}{d\boldsymbol{\theta}} F(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) &= \frac{d}{d\boldsymbol{\theta}} \frac{d}{d\beta} F(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta), \\ \frac{d}{d\beta} \frac{\partial}{\partial \boldsymbol{\theta}} F(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) &= \frac{d}{d\boldsymbol{\theta}} \frac{\partial}{\partial \beta} F(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta), \quad \text{by Eq. (13).} \end{aligned}$$

By then applying this equality in  $\beta = 0$  and  $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ , we obtain the EP gradient formula (Eq. (2)) for complex differentiation:

$$\left. \frac{d}{d\beta} \right|_{\beta=0} \left( \frac{\partial F}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) \right) = \frac{d}{d\boldsymbol{\theta}} \frac{\partial F}{\partial \beta}(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta) = \frac{d\mathcal{L}}{d\boldsymbol{\theta}},$$

which concludes the proof.  $\square$

### A.3 Proof of Theorem 1

**Theorem 1** (Exact gradient from finite teaching signals). *Assuming that the conditions of Lemma 1 are met and let  $|\beta| > 0$  be the radius of a circular path around 0 in  $\mathbb{C}$  contained in the open set  $U$  on which the fixed point  $\mathbf{s}_{\boldsymbol{\theta}, \beta}^*$  is defined. Further assume that this path is parameterized by  $t \in [0, T] \mapsto \beta(t) = |\beta|e^{2i\pi t/T}$ , where  $i$  is the imaginary unit. Then the loss gradient is given by:*

$$\frac{d\mathcal{L}}{d\boldsymbol{\theta}} = \frac{1}{T|\beta|} \int_0^T \frac{\partial F}{\partial \boldsymbol{\theta}}(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta(t)}^*, \beta(t)) e^{-2i\pi t/T} dt. \quad (14)$$

*Proof.* By assumption the fixed point  $\beta \mapsto \mathbf{s}_{\boldsymbol{\theta}, \beta}^*$  is defined on an open set  $U$  (by the holomorphic implicit function theorem) containing the disk of radius  $|\beta|$  centered around 0. In particular, the function  $\beta \in U \mapsto \partial_{\boldsymbol{\theta}} F(\boldsymbol{\theta}, \mathbf{s}_{\boldsymbol{\theta}, \beta}^*, \beta)$ , is also holomorphic by composition. The left hand side of

Eq. (4) can thus be computed with the Cauchy formulas (Eq. (12) with  $f = \partial_{\theta} F$ ,  $n = 1$ ,  $a = 0$ ), and  $\gamma$  an arbitrary closed path leading around zero once and counterclockwise in  $U$ :

$$\left. \frac{d}{d\beta} \right|_{\beta=0} \left( \frac{\partial F}{\partial \theta}(\theta, s_{\theta, \beta}^*, \beta) \right) = \frac{1}{2i\pi} \oint_{\gamma} \frac{1}{\beta^2} \frac{\partial F}{\partial \theta}(\theta, s_{\theta, \beta}^*, \beta) d\beta \quad (15)$$

To obtain Eq. (14), we choose  $\gamma$  as a circular path in the complex plane with radius  $|\beta| > 0$  parameterized by time  $t \in [0, T] \mapsto \beta(t) = |\beta|e^{2i\pi t/T}$ , where  $T$  is a full period. After the change of variable  $d\beta = (2i\pi\beta(t)/T)dt$  in Eq. (15), and using Lemma 1, the loss gradient is given by:

$$\begin{aligned} \frac{d\mathcal{L}}{d\theta} &= \frac{1}{2i\pi} \oint_{\gamma} \frac{1}{\beta^2} \frac{\partial F}{\partial \theta}(\theta, s_{\theta, \beta}^*, \beta) d\beta \\ &= \frac{1}{2i\pi} \int_0^T \frac{1}{\beta(t)^2} \frac{\partial F}{\partial \theta}(\theta, s_{\theta, \beta(t)}^*, \beta(t)) \left( \frac{2i\pi\beta(t)}{T} \right) dt \\ &= \frac{1}{T} \int_0^T \frac{1}{\beta(t)} \frac{\partial F}{\partial \theta}(\theta, s_{\theta, \beta(t)}^*, \beta(t)) dt \\ &= \frac{1}{T|\beta|} \int_0^T \frac{\partial F}{\partial \theta}(\theta, s_{\theta, \beta(t)}^*, \beta(t)) e^{-2i\pi t/T} dt. \end{aligned}$$

□

#### A.4 Roles of real and imaginary parts in the learning rule

Recall that for the continuous Hopfield network case the partial derivative of  $F$  with respect to a parameter  $w_{ij}$  is the product of pre and post activation (Eq. (6)), so that applying Eq. (4) yields:

$$\left. \frac{d\mathcal{L}}{dw_{ij}} \right|_{\beta=0} = \left. \frac{d}{d\beta} \right|_{\beta=0} \underbrace{\left( \frac{\partial F}{\partial w_{ij}}(\theta, s_{\beta}^*, \beta) \right)}_{=-\sigma(s_{i, \beta}^*)\sigma(s_{j, \beta}^*)} = - \left. \frac{d(\sigma(s_{i, \beta}^*)\sigma(s_{j, \beta}^*))}{d\beta} \right|_{\beta=0},$$

which can further be expressed as:

$$\left. \frac{d(\sigma(s_{i, \beta}^*)\sigma(s_{j, \beta}^*))}{d\beta} \right|_{\beta=0} = \left( \sigma(s_{i, \beta}^*) \frac{d\sigma(s_{j, \beta}^*)}{d\beta} \right) \Big|_{\beta=0} + \left( \sigma(s_{j, \beta}^*) \frac{d\sigma(s_{i, \beta}^*)}{d\beta} \right) \Big|_{\beta=0}. \quad (16)$$

Using the same assumptions as Section 3, the map  $\beta \in U \mapsto s_{i, \beta}^*$  is holomorphic, and so is the map  $\beta \in U \mapsto \sigma(s_{i, \beta}^*)$  by composition. We can thus expand it in a power series around zero:

$$\sigma(s_{i, \beta}^*) = \sum_{k=0}^{\infty} \frac{\beta^k}{k!} \left. \frac{d^k \sigma(s_{i, \beta}^*)}{d\beta^k} \right|_{\beta=0}.$$

We can then separate the sum into the real and imaginary parts because the series converge absolutely.

Assuming that  $\beta = |\beta|e^{2i\pi t/T}$ , and applying the Euler formula, we obtain:

$$\begin{aligned} \operatorname{Re}(\sigma(s_{i, \beta}^*)) &= \sum_{k=0}^{\infty} \cos\left(\frac{2k\pi t}{T}\right) \frac{|\beta|^k}{k!} \left. \frac{d^k \sigma(s_{i, \beta}^*)}{d\beta^k} \right|_{\beta=0}, \\ \operatorname{Im}(\sigma(s_{i, \beta}^*)) &= \sum_{k=1}^{\infty} \sin\left(\frac{2k\pi t}{T}\right) \frac{|\beta|^k}{k!} \left. \frac{d^k \sigma(s_{i, \beta}^*)}{d\beta^k} \right|_{\beta=0}. \end{aligned} \quad (17)$$

Therefore, the first derivative ( $k = 1$ ) with respect to  $\beta$  in Eq. (16) can be obtained by either projecting the real part against the cosine function, or imaginary part against the sine function:

$$\begin{aligned} \left. \frac{d\sigma(s_{i, \beta}^*)}{d\beta} \right|_{\beta=0} &= \frac{2}{|\beta|T} \int_0^T \operatorname{Re}(\sigma(s_{i, \beta}^*)) \cos\left(\frac{2\pi t}{T}\right) dt, \\ &= \frac{2}{|\beta|T} \int_0^T \operatorname{Im}(\sigma(s_{i, \beta}^*)) \sin\left(\frac{2\pi t}{T}\right) dt, \end{aligned}$$

by orthogonality of the family  $((t \mapsto \cos(\frac{2k\pi t}{T}))_{k \geq 0}, (t \mapsto \sin(\frac{2k\pi t}{T}))_{k \geq 0})$  in  $L^2[0, T]$ . Note that the higher order derivatives with respect to  $\beta$  can be obtained as well by projecting against the corresponding cosine or sine function. The same holds for index  $j$  by symmetry. As an interesting final note, if we define  $\text{Re}_1(\sigma(s_{i,\beta}^*))$  and  $\text{Im}_1(\sigma(s_{i,\beta}^*))$ , the first order contributions in  $\beta$  to the real and imaginary parts of the neural activity, we find that they are the only ones to contribute to the gradient computation. We can appreciate that they are related through  $\text{Im}_1(\sigma(s_{i,\beta}^*)) = -\frac{T}{2\pi} \frac{d}{dt} \text{Re}_1(\sigma(s_{i,\beta}^*))$ , where the time derivative is at the scale of the teaching signal.

## A.5 Derivation of the bias term

Recall the definition of  $\beta_k := |\beta| e^{2i\pi k/N}$ , for  $k \in [0, \dots, N-1]$ ,  $N \geq 2$ , and the gradient estimate (Eq. (8)):

$$\hat{\nabla}(N) := \frac{1}{N|\beta|} \sum_{k=0}^{N-1} \frac{\partial F}{\partial \theta}(\theta, \mathbf{s}_{\beta_k}^*, \beta_k) e^{-2i\pi k/N}.$$

For simplicity of notation, we rewrite  $\partial_{\theta} F(\beta) := \frac{\partial F}{\partial \theta}(\theta, \mathbf{s}_{\beta}^*, \beta)$ . The function  $\beta \mapsto \partial_{\theta} F(\beta)$  is holomorphic on an open set  $U$  including zero, and so is  $\beta \mapsto \beta_{\beta}^*$  by the holomorphic implicit function theorem. We assume the  $\beta_k$  are included in  $U$ , so that we can expand  $\partial_{\theta} F(\beta_k)$  in a power series around zero:

$$\partial_{\theta} F(\beta_k) = \sum_{p=0}^{\infty} \frac{\beta_k^p}{p!} \left[ \frac{d^p}{d\beta^p} \partial_{\theta} F \right](0),$$

we define  $C_p := \left[ \frac{d^p}{d\beta^p} \partial_{\theta} F \right](0)$ . The quantity of interest is  $C_1$ , since it is the gradient of the loss (Eq. (4))

$$\begin{aligned} \partial_{\theta} F(\beta_k) &= C_0 + \beta_k C_1 + \sum_{p=2}^{\infty} \frac{\beta_k^p}{p!} C_p \\ \frac{\partial_{\theta} F(\beta_k)}{\beta_k} &= C_0 \beta_k^{-1} + C_1 + \sum_{p=2}^{\infty} \frac{\beta_k^{p-1}}{p!} C_p \\ \frac{1}{N} \sum_{k=0}^{N-1} \frac{\partial_{\theta} F(\beta_k)}{\beta_k} &= C_1 + C_0 \frac{1}{N} \sum_{k=0}^{N-1} \beta_k^{-1} + \frac{1}{N} \sum_{k=0}^{N-1} \sum_{p=2}^{\infty} \frac{\beta_k^{p-1}}{p!} C_p. \end{aligned}$$

The sum symbols on the right can be interchanged thanks to the absolute convergence of the power series.

$$\begin{aligned} \frac{1}{N} \sum_{k=0}^{N-1} \frac{\partial_{\theta} F(\beta_k)}{\beta_k} &= C_1 + C_0 \frac{1}{N} \sum_{k=0}^{N-1} \beta_k^{-1} + \sum_{p=2}^{\infty} \frac{C_p}{p!} \frac{1}{N} \sum_{k=0}^{N-1} \beta_k^{p-1} \\ \frac{1}{N} \sum_{k=0}^{N-1} \frac{\partial_{\theta} F(\beta_k)}{\beta_k} &= C_1 + C_0 \frac{1}{N|\beta|} \sum_{k=0}^{N-1} e^{-2i\pi k/N} + \sum_{p=1}^{\infty} \frac{C_{p+1}}{(p+1)!} \frac{|\beta|^p}{N} \sum_{k=0}^{N-1} e^{2i\pi p k/N}. \end{aligned}$$

It remains to evaluate the geometric sums of the form  $\sum_{k=0}^{N-1} e^{2i\pi p k/N}$  for  $p = -1$  and  $p \geq 1$ . If  $N$  divides  $p$ , i.e  $p \equiv 0 \pmod{N}$ , then we can write  $p = Nq$  and we have:

$$\sum_{k=0}^{N-1} e^{2i\pi p k/N} = \sum_{k=0}^{N-1} e^{2i\pi q N k/N} = \sum_{k=0}^{N-1} e^{2i\pi q k} = \sum_{k=0}^{N-1} 1 = N.$$

If  $N$  does not divide  $p$ , then the geometric sum of ratio  $e^{2i\pi p/N}$  can be computed:

$$\sum_{k=0}^{N-1} e^{2i\pi p k/N} = \frac{1 - (e^{2i\pi p/N})^N}{1 - e^{2i\pi p/N}} = \frac{1 - e^{2i\pi p}}{1 - e^{2i\pi p/N}} = \frac{1 - 1}{1 - e^{2i\pi p/N}} = 0.$$

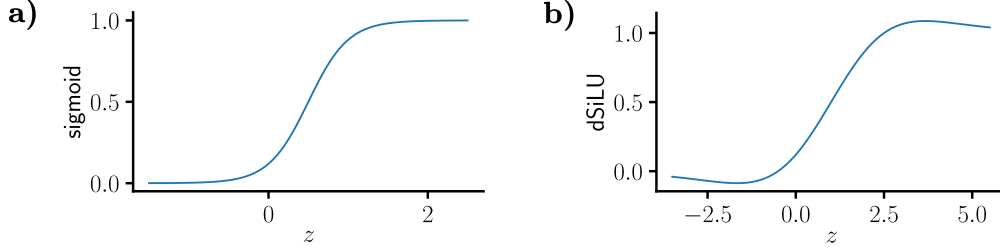


Figure 5: **a)** The shifted sigmoid we used in multi-layered perceptrons experiments. **b)** The dSiLU we used in CNNs experiments.

660 We thus have that:

$$\begin{aligned} \frac{1}{N} \sum_{k=0}^{N-1} \frac{\partial_{\theta} F(\beta_k)}{\beta_k} &= C_1 + C_0 \frac{1}{N|\beta|} \underbrace{\sum_{k=0}^{N-1} e^{-2i\pi k/N}}_{=0} + \sum_{p=1}^{\infty} \frac{C_{p+1}}{(p+1)!} \frac{|\beta|^p}{N} \underbrace{\sum_{k=0}^{N-1} e^{2i\pi p k/N}}_{=0 \text{ when } p \neq 0 \text{ (N)}} \\ \frac{1}{N} \sum_{k=0}^{N-1} \frac{\partial_{\theta} F(\beta_k)}{\beta_k} &= C_1 + \sum_{p \equiv 0 \text{ (N)}}^{\infty} \frac{C_{p+1} |\beta|^p}{(p+1)!}, \end{aligned}$$

661 which is the result of Eq. (9).

## 662 A.6 Derivation of the online estimate

663 Recall the formula of the online estimate (Eq. (10)):

$$\tilde{\nabla}(T_{\text{plas}}) := -\frac{1}{T_{\text{plas}}|\beta|} \int_0^{T_{\text{plas}}} \sigma_i(t) \sigma_j(t) e^{-2i\pi t/T_{\text{osc}}} dt.$$

664 If  $T_{\text{dyn}} \ll T_{\text{osc}}$ , the product of activities can be replaced by its value at the fixed point, and an exact  
 665 gradient is computed after each period (Eq. (7)). Then if  $T_{\text{osc}} \ll T_{\text{plas}}$ , the integral can be divided  
 666 into an integer amount of completed periods plus a remainder:  $T_{\text{plas}} = kT_{\text{osc}} + T_{\text{rem}}$ , where  $k \in \mathbb{N}$   
 667 and  $T_{\text{rem}} < T_{\text{osc}}$ . We then have by periodicity that:

$$\tilde{\nabla}(T_{\text{plas}}) = \underbrace{\frac{kT_{\text{osc}}}{kT_{\text{osc}} + T_{\text{rem}}}}_{\rightarrow 1 \text{ when } T_{\text{plas}} \rightarrow \infty} \frac{d\mathcal{L}}{dw_{ij}} - \underbrace{\frac{1}{kT_{\text{osc}} + T_{\text{rem}}}}_{\rightarrow 0 \text{ when } T_{\text{plas}} \rightarrow \infty} \frac{1}{|\beta|} \int_0^{T_{\text{rem}}} \sigma_i^*(t) \sigma_j^*(t) e^{-2i\pi t/T_{\text{osc}}} dt.$$

668 In this way, when averaging over large  $T_{\text{plas}}$ , the number of completed cycles outweighs the current  
 669 period. Thus, by simply averaging over many oscillation cycles allows estimating gradients without  
 670 explicit separate phases.

## 671 B Detailed architecture

### 672 B.1 Dynamics for multi-layer perceptrons

673 Assuming a number of  $L$  layers, we note  $\mathbf{s}_l$  the subset of units in layer  $l$ , with  $\mathbf{s}_0 = \mathbf{x}$  and  $\mathbf{y}$  the one  
 674 hot class label. We note  $\mathbf{W}_l$ , and  $\mathbf{b}_l$  the weight and biases of layer  $l \geq 1$ . The energy function  $F$  for  
 675 a MLP optimizing the cross entropy loss reads:

$$F(\theta, \mathbf{s}, \beta, \mathbf{y}) = \sum_{l=1}^{L-1} \left[ \frac{1}{2} \|\mathbf{s}_l\|^2 - \sigma(\mathbf{s}_{l-1})^\top \cdot \mathbf{W}_l \cdot \sigma(\mathbf{s}_l) - \mathbf{b}_l^\top \cdot \sigma(\mathbf{s}_l) \right] - \beta \mathbf{y}^\top \cdot \log(\mathbf{s}_L).$$

676 The activation function used for multi-layer perceptrons is the shifted sigmoid  $z \mapsto 1/(1 + e^{-4z+2})$   
 677 (Fig. 5a). We use the layer-wise discrete dynamics introduced by [S3, S4], which read:

$$\begin{cases} \mathbf{s}_l & \leftarrow \sigma \left( \mathbf{W}_l \mathbf{s}_{l-1} + \mathbf{W}_{l+1}^\top \mathbf{s}_{l+1} + \mathbf{b}_l + \boldsymbol{\eta}_l \right), & \text{for } 1 \leq l \leq L-2 \\ \mathbf{s}_{L-1} & \leftarrow \sigma \left( \mathbf{W}_{L-1} \mathbf{s}_{L-2} + \beta \mathbf{W}_L^\top (\mathbf{y} - \mathbf{s}_L) + \mathbf{b}_{L-1} + \boldsymbol{\eta}_{L-1} \right), \\ \mathbf{s}_L & \leftarrow \text{Softmax}(\mathbf{W}_L \mathbf{s}_{L-1} + \mathbf{b}_L), \end{cases}$$

678 where  $\boldsymbol{\eta}_l$  is an optional Gaussian noise added for Fig. 3c and Table 1. The noise was sampled at each  
 679 time step.

## 680 B.2 Dynamics for convolutional neural networks

681 The activation function used for CNNs is a sigmoid-weighted linear unit [S5] (Fig. 5b):

$$\text{dSiLU}(z) := \left(\frac{z}{2}\right) \frac{1}{1 + e^{-z}} + \left(1 - \frac{z}{2}\right) \frac{1}{1 + e^{-z+2}}.$$

682 We denote by  $\mathcal{P}$  the pooling operation, and  $\tilde{\mathcal{P}}$  the corresponding unpooling operation. ‘ $*$ ’ denotes the  
 683 convolution when preceded by  $\mathbf{W}$  and transpose convolution when preceded by  $\mathbf{W}^\top$ . The energy  
 684 function  $F$  for a CNN optimizing the cross entropy loss reads:

$$\begin{aligned} F(\boldsymbol{\theta}, \mathbf{s}, \beta, \mathbf{y}) = & \sum_{l \in \{\text{Conv}\}} \left[ \frac{1}{2} \|\mathbf{s}_l\|^2 - \sigma(\mathbf{s}_{l-1})^\top \cdot \mathcal{P}(\mathbf{W}_l * \sigma(\mathbf{s}_l)) - \mathbf{b}_l^\top \cdot \sigma(\mathbf{s}_l) \right] \\ & \sum_{l \in \{\text{FC}\}} \left[ \frac{1}{2} \|\mathbf{s}_l\|^2 - \sigma(\mathbf{s}_{l-1})^\top \cdot \mathbf{W}_l \cdot \sigma(\mathbf{s}_l) - \mathbf{b}_l^\top \cdot \sigma(\mathbf{s}_l) \right] - \beta \mathbf{y}^\top \cdot \log(\mathbf{s}_L). \end{aligned}$$

685 We use the layer-wise discrete dynamics introduced by [S3, S4], which read:

$$\begin{cases} \mathbf{s}_l & \leftarrow \sigma \left( \mathcal{P}(\mathbf{W}_l * \mathbf{s}_{l-1}) + \mathbf{W}_{l+1}^\top * \tilde{\mathcal{P}}(\mathbf{s}_{l+1}) + \mathbf{b}_l \right), & \text{for } l \in \{\text{Conv layers}\} \\ \mathbf{s}_l & \leftarrow \sigma \left( \mathbf{W}_l \mathbf{s}_{l-1} + \mathbf{W}_{l+1}^\top \mathbf{s}_{l+1} + \mathbf{b}_l \right), & \text{for } l \in \{\text{FC layers}\} \\ \mathbf{s}_{L-1} & \leftarrow \sigma \left( \mathbf{W}_{L-1} \mathbf{s}_{L-2} + \beta \mathbf{W}_L^\top (\mathbf{y} - \mathbf{s}_L) + \mathbf{b}_{L-1} \right), \\ \mathbf{s}_L & \leftarrow \text{Softmax}(\mathbf{W}_L \mathbf{s}_{L-1} + \mathbf{b}_L). \end{cases}$$

686 We used Softmax pooling [S6] with a tunable temperature  $\tau$ , instead of the non-holomorphic Max  
 687 pooling. The output  $y$  of Softmax pooling of an input  $\mathbf{x}$  is defined for a kernel neighbourhood  $\mathbf{R}$  by:

$$y = \sum_{i \in \mathbf{R}} \left( \frac{e^{x_i/\tau}}{\sum_{j \in \mathbf{R}} e^{x_j/\tau}} \right) x_i.$$

688 Note that Softmax pooling interpolates between Average pooling ( $\tau \rightarrow \infty$ ) and Max pooling ( $\tau \rightarrow 0$ ).

## 689 C Layer-wise comparison of the gradient in a deep network

690 Here we show in Fig. 6 the complete layer-wise cosine similarities between the estimates of holomor-  
 691 phic EP for various  $N$  and the true gradient computed by automatic differentiation.

## 692 D Dynamical stability in the complex plane

693 We show in Fig. 7 how the area in  $\mathbb{C}$  where the fixed point empirically exists varies with different  
 694 architecture choices. The data used for each panel is a digit from the MNIST dataset. As in Fig.2b),  
 695 dark blue means that the fixed point exists, whereas light areas denote divergence. These diverging  
 696 areas could be due to the poles of the activation functions used. For example, the sigmoid function

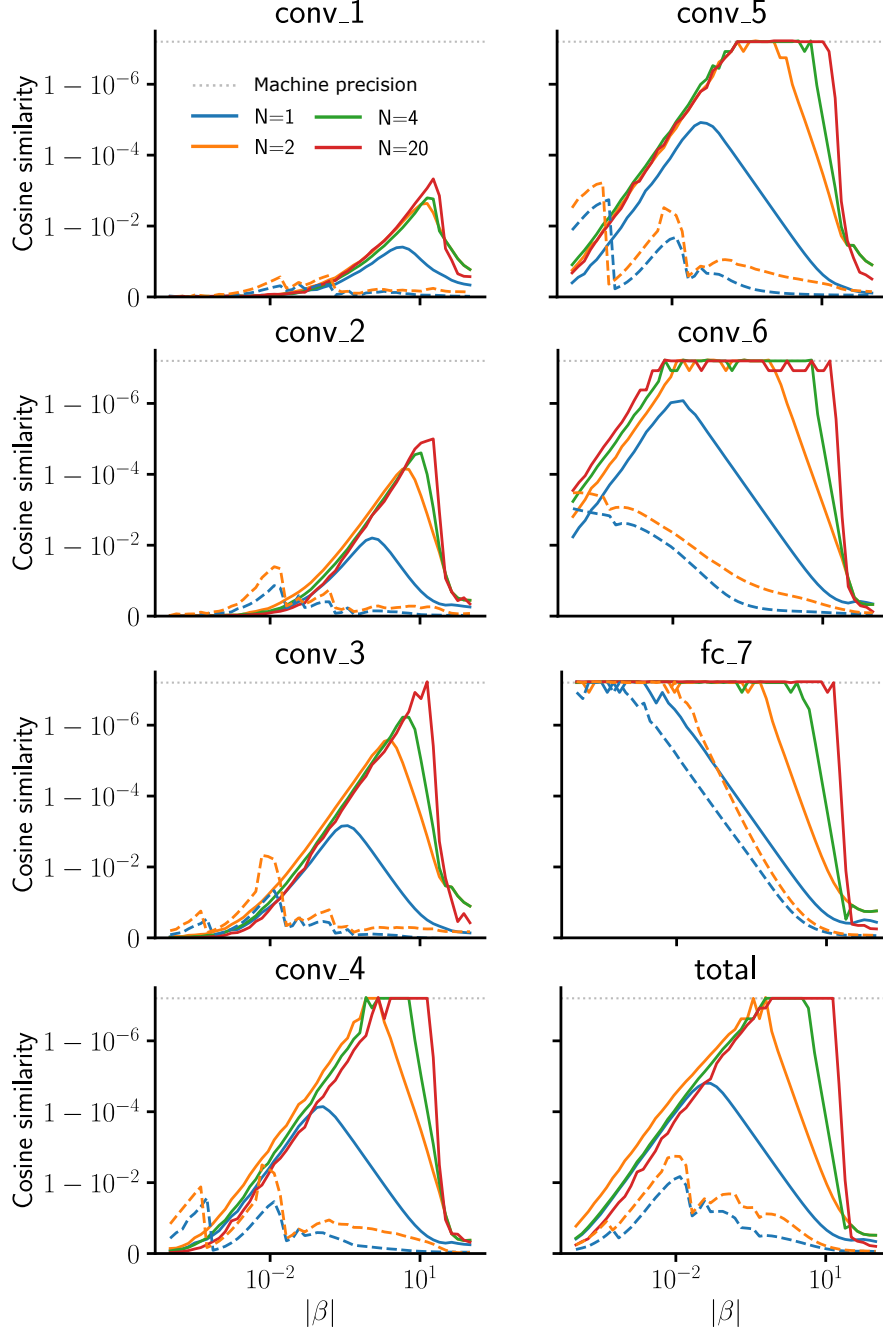


Figure 6: The complete layer-wise cosine similarity of Fig. 4a).

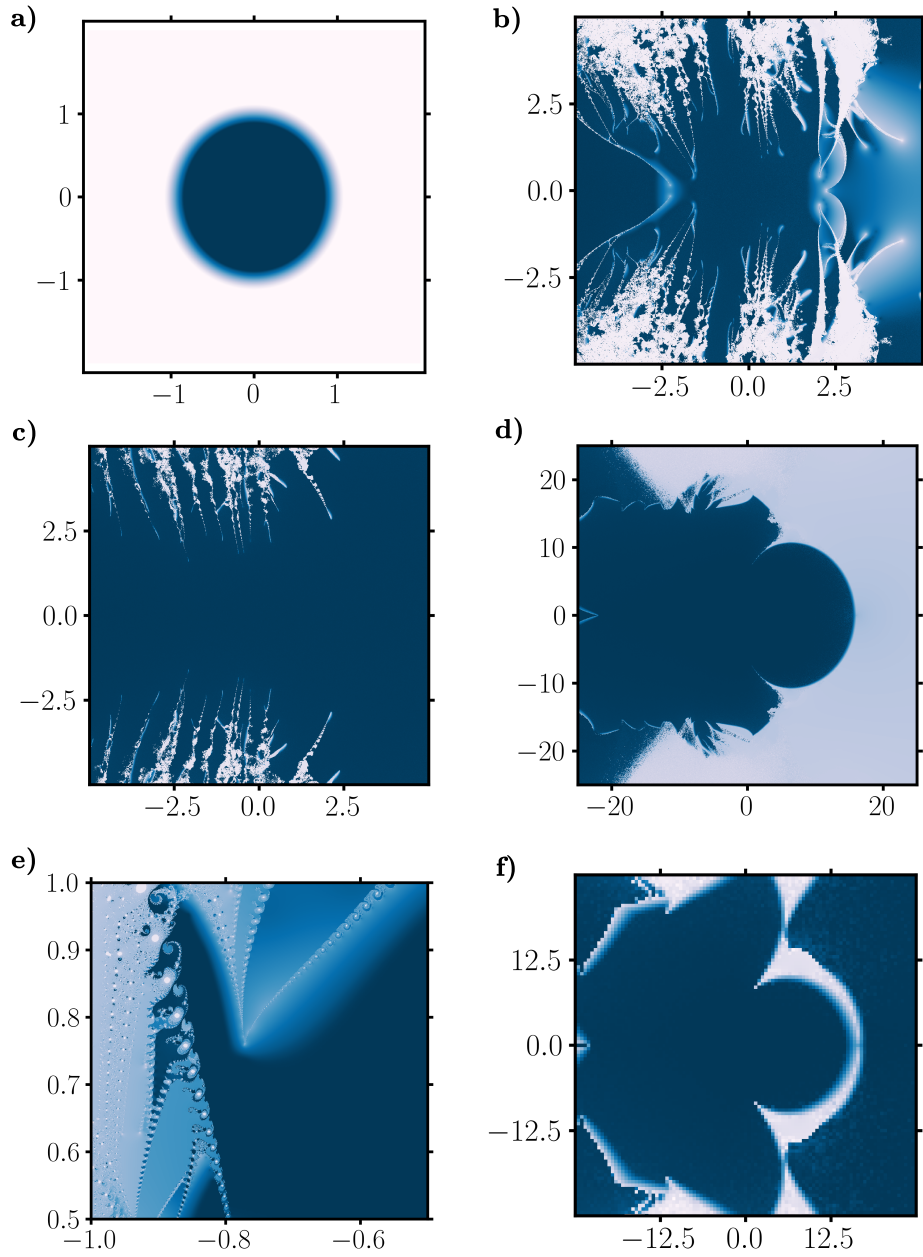


Figure 7: Map of convergence to a fixed point for complex  $\beta$  in various settings. **a)** MLP with linear activation function and low weight initialization. **b)** MLP with shifted sigmoid activation and default weight initialization. **c)** Same as b) but with reduced weight initialization. **d)** Same as b) but with dSiLU activation function. **e)** Zoom at a frontier between stable and unstable regions. **f)** Small CNN with dSiLU activation and Softmax pooling.



697  $z \mapsto 1/(1 + e^{-z})$  has  $\{(2k + 1)i\pi; k \in \mathbb{Z}\}$  as a set of poles where it diverges. Although we did not  
698 systematically study this phenomenon in this work, we strongly suspect that these unstable areas are  
699 partly the result of the teaching signal being too strong or the weights being poorly distributed, thereby  
700 driving the complex neural activities near to the poles. To some extent, the poles can be brought  
701 farther by introducing a coefficient in the exponential, but it results in flatter activation functions  
702 on the real axis, so a trade-off should be found. In practice, we found that choosing reasonably the  
703 activation function, weight initialization, and teaching radius  $|\beta|$  lead to enough stable areas around 0  
704 to compute the gradient.

## 705 E Hyperparameters

### 706 E.1 MNIST experiments

707 The digits were rescaled by 255 and flattened. The hyperparameters used for training are reported in  
708 Table 3, and the training errors are reported in Table 4.

Table 3: Hyperparameters used for the MNIST training experiment of Table 1.

Hyperparameter	Classic EP	hEP	Online hEP
Batch size	20	20	20
Learning rate	5e-2	5e-2	5e-2
Epochs	50	50	50
$ \beta $	0.1 and 0.4	0.4	0.4
$T_{\text{free}}$	350	200	200*
$T_{\text{nudge}}$	350	50	N/A
$T_{\text{osc}}$	N/A	N/A	300
$T_{\text{plas}}$	N/A	N/A	900
$N$	N/A	10	10
Noise**	4e-2	4e-2	4e-2

\* Only used for evaluation

\*\* Standard deviation of the Gaussian noise for experiments with noise.

Table 4: MNIST training and validation errors for classic EP [10], hEP, and online hEP, with and without noise. All results are averages ( $n = 3$ )  $\pm$  one standard deviation.

Noise	Class. EP, $ \beta  = 0.1$		Class. EP, $ \beta  = 0.4$		hEP, $ \beta  = 0.4$		Online hEP	
	Train (%)	Val (%)	Train (%)	Val (%)	Train (%)	Val (%)	Train (%)	Val. (%)
No	0.05 $\pm 0.02$	1.87 $\pm 0.01$	0.19 $\pm 0.05$	2.24 $\pm 0.05$	0.02 $\pm 0.01$	1.97 $\pm 0.08$	0.11 $\pm 0.01$	2.05 $\pm 0.02$
Yes	88.8 $\pm 0.0$	88.7 $\pm 0.0$	1.96 $\pm 0.2$	3.01 $\pm 0.1$	0.14 $\pm 0.03$	1.96 $\pm 0.07$	0.13 $\pm 0.03$	1.91 $\pm 0.16$

### 709 E.2 Large-scale experiments

710 In the training experiments for CIFAR-10, CIFAR-100, and ImageNet  $32 \times 32$ , the training data was  
711 normalized, then augmented with 50% chance random horizontal flips, resized to  $36 \times 36$  resolution  
712 with padding, and cropped randomly back to  $32 \times 32$ . The optimizer used was stochastic gradient  
713 descent with momentum. Pooling was applied at all layers for the five-layer CNN, and every other  
714 layer starting with the first layer in the seven-layer CNN.

## 715 F Simulations details

716 All simulations were performed on an in-house GPU cluster or workstations. Each simulation in  
717 Table 2 was run in parallel on four NVIDIA V100 GPUs. The training runs on ImageNet  $32 \times 32$   
718 took 5.5 days each for EP, and a few hours for BP. The runs on CIFAR-10 and CIFAR-100 took one  
719 day on average depending on the architecture (5 or 7 layers) and the number of time steps used for

Table 5: Hyperparameters used for the VGG training experiments of Table 2 and Fig.4c.

Hyperparameter	CIFAR-10	CIFAR-100	ImageNet $32 \times 32$	CIFAR-10 (Fig.4c)
Batch size	128	128	256	128
Channel sizes		[128, 256, 512, 512]		[128, 128, 256, 256, 512, 512]
Kernel sizes		[3, 3, 3, 3]		[3, 3, 3, 3, 3, 3]
Strides		[1, 1, 1, 1]		[1, 1, 1, 1, 1, 1]
Paddings		[1, 1, 1, 0]		[1, 1, 1, 0, 1, 0]
SoftPool window		$2 \times 2$		$2 \times 2$
SoftPool stride		2		2
SoftPool temp.		1		10
Initial LRs*		$[25, 15, 10, 8, 5] \times 1e-2$		$[5, 4, 4, 3, 3, 2, 2] \times 1e-2$
Final LRs		$[25, 15, 10, 8, 5] \times 1e-9$		$[5, 4, 4, 3, 3, 2, 2] \times 1e-9$
Weight decay	2e-3	1e-2	5e-4	$[5, 5, 5, 5, 5, 5, 10] \times 1e-4$
Momentum	0.9	0.9	0.9	0.9
Epochs	90	90	90	90
$ \beta $	1.0	1.0	1.0	1.0
$T_{\text{free}}$	250	250	250	260
$T_{\text{nudge}}$	60	60	60	60
$N$	2	2	2	2 and 4

\* Learning rates were decayed with cosine annealing without restart [S7].

the dynamics. The use of complex numbers, although seamlessly implementable with Jax, results in longer simulation times due to the 64 bit-precision requirement (32 bit-precision for real and imaginary parts respectively).

## 723 **Supplementary References**

- 724 [S1] Walter Appel. Mathematics for physics and physicists. 2007.
- 725 [S2] Henri Cartan. *Théorie élémentaire des fonctions analytiques d'une ou plusieurs variables*  
726 *complexes: Avec le concours de Reiji Takahashi*. Hermann, 1961.
- 727 [S3] Maxence Ernoult, Julie Grollier, Damien Querlioz, Yoshua Bengio, and Benjamin Scellier.  
728 Updates of equilibrium prop match gradients of backprop through time in an rnn with static input.  
729 *Advances in neural information processing systems*, 32, 2019.
- 730 [S4] Axel Laborieux, Maxence Ernoult, Benjamin Scellier, Yoshua Bengio, Julie Grollier, and Damien  
731 Querlioz. Scaling equilibrium propagation to deep convnets by drastically reducing its gradient  
732 estimator bias. *Frontiers in neuroscience*, 15:129, 2021.
- 733 [S5] Stefan Elfving, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network  
734 function approximation in reinforcement learning. *Neural Networks*, 107:3–11, 2018.
- 735 [S6] A Stergiou, R Poppe, and G Kalliatakis. Refining activation downsampling with softpool. arxiv  
736 2021. *arXiv preprint arXiv:2101.00440*.
- 737 [S7] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*  
738 *arXiv:1711.05101*, 2017.