

Scaling EqProp to Deep ConvNets by reducing its gradient estimator bias

Axel Laborieux*, Maxence Ernoult*, Benjamin Scellier, Yoshua Bengio, Julie Grollier, Damien Querlioz

* Corresponding authors: {axel.laborieux, maxence.ernoult}@universite-paris-saclay.fr

Summary

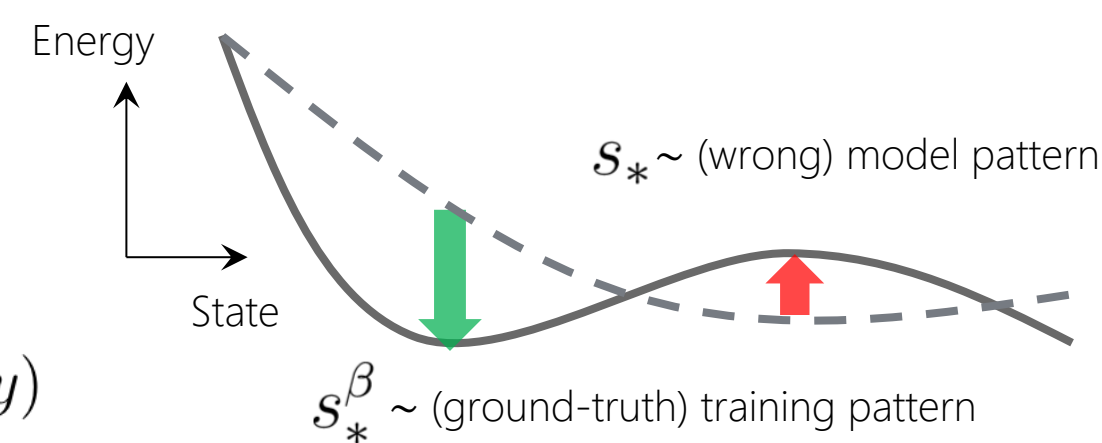
- Equilibrium Propagation (EP) is a biologically-inspired counterpart of Backpropagation Through Time (BPTT). However :
 1. Only shown to work on MNIST
 2. Optimizes the squared error loss
- **This work** : identify and cancel a bias in the gradient estimate
 1. Scale to deeper architectures and closely match BPTT on CIFAR10
 2. Adapt the model to optimize the cross-entropy loss
 3. Still works with untied weights provided an alignment mechanism

Equilibrium Propagation [1]

$$s_{t+1} = \frac{\partial \Phi}{\partial s}(x, s_t, \theta)$$

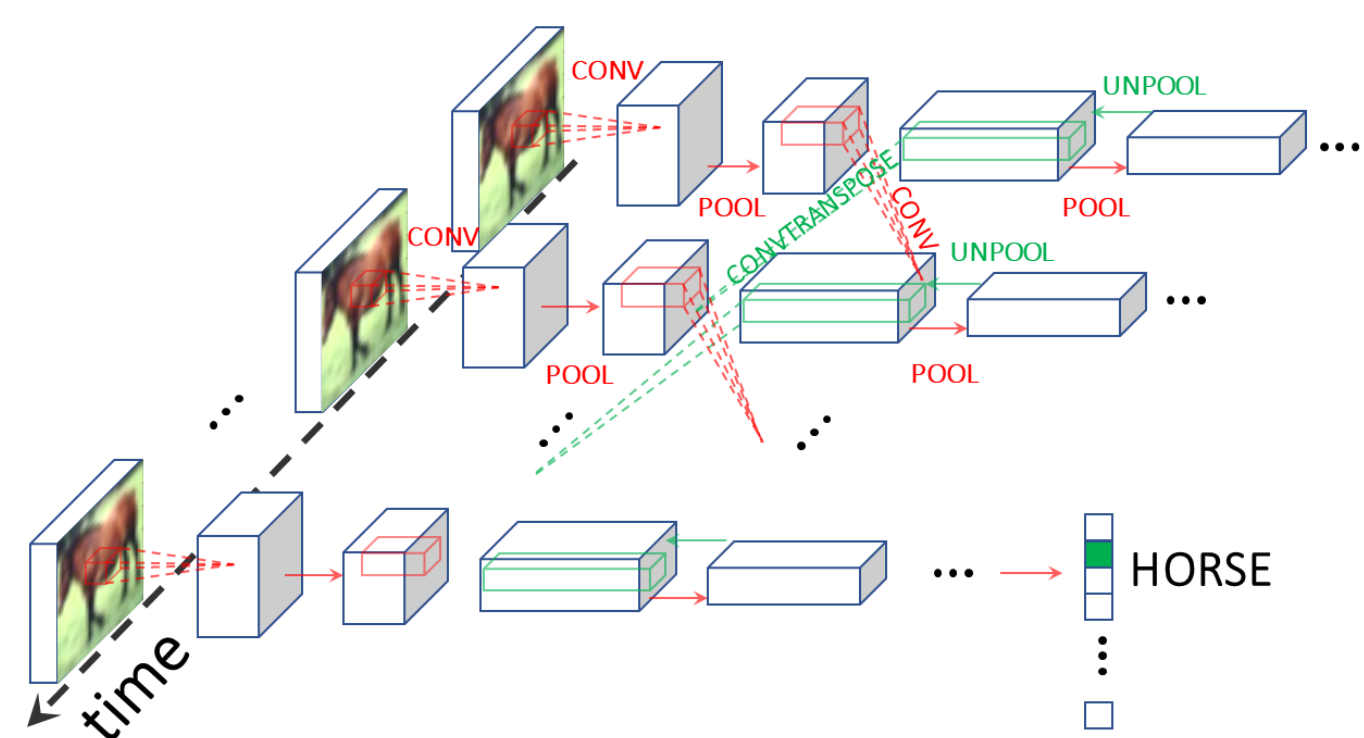
$$\text{until } s_* = \frac{\partial \Phi}{\partial s}(x, s_*, \theta)$$

$$s_{t+1}^\beta = \frac{\partial \Phi}{\partial s}(x, s_t^\beta, \theta) - \beta \frac{\partial \ell}{\partial s}(s_t^\beta, y)$$



$$\widehat{\nabla}^{\text{EP}}(\beta) \triangleq \frac{1}{\beta} \left(\frac{\partial \Phi}{\partial \theta}(x, s_*^\beta, \theta) - \frac{\partial \Phi}{\partial \theta}(x, s_*, \theta) \right)$$

Convolutional Architecture [2]



$$\begin{cases} s_{t+1}^n = \sigma(\mathcal{P}(w_{n-1} * s_t^{n-1}) + \tilde{w}_{n+1} * \mathcal{P}^{-1}(s_t^{n+1})), & \text{for convolutional layers,} \\ s_{t+1}^n = \sigma(w_{n-1} \cdot s_t^{n-1} + w_{n+1}^\top \cdot s_t^{n+1}), & \text{for fully connected layers.} \end{cases}$$



<https://arxiv.org/pdf/2006.03824.pdf>

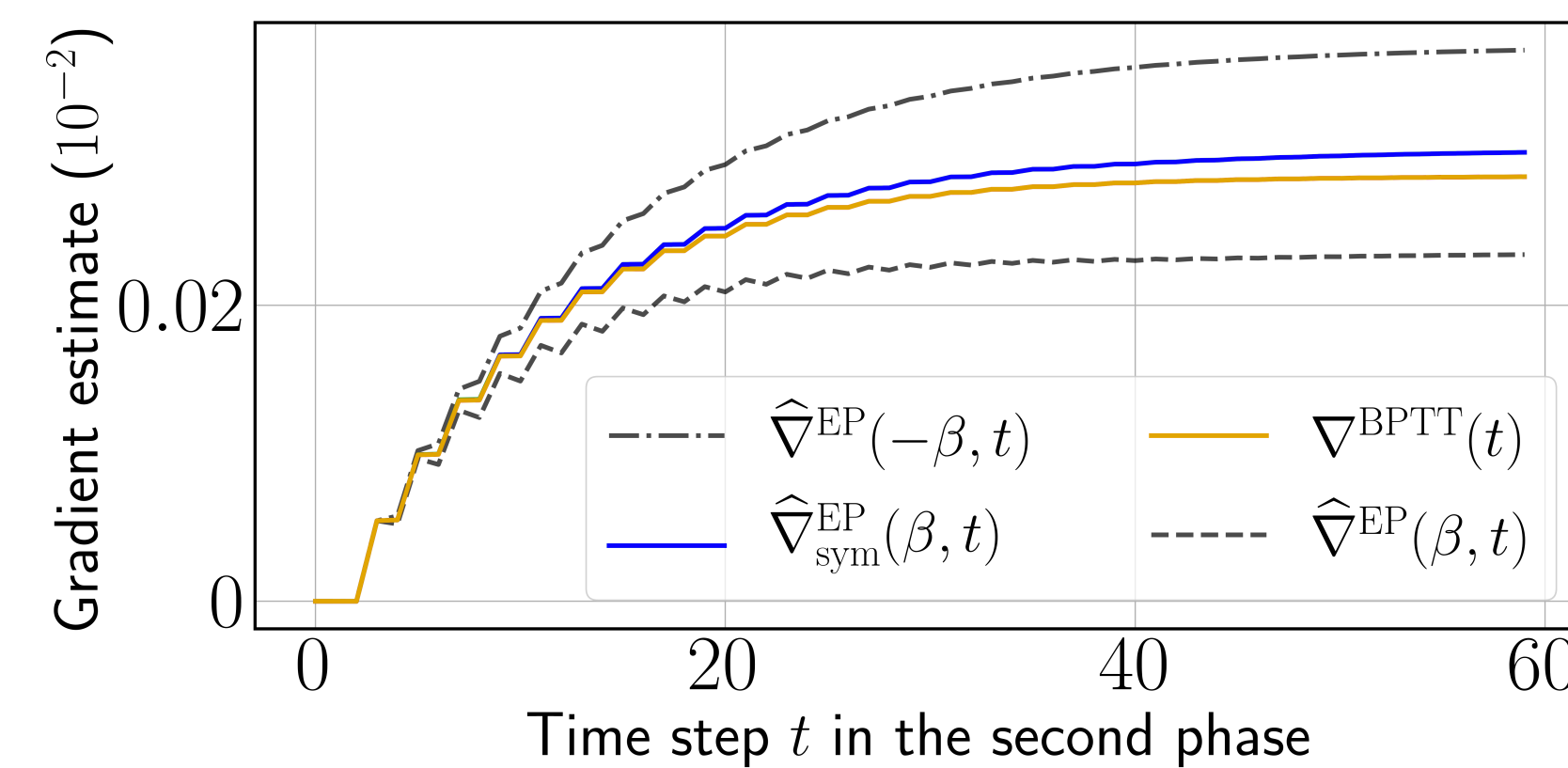
<https://github.com/Laborieux-Axel/Equilibrium-Propagation>

Symmetric estimate of the loss gradient

$$\widehat{\nabla}^{\text{EP}}(\beta) = -\frac{\partial \mathcal{L}^*}{\partial \theta} + O(\beta) : \text{first order bias in } \beta$$

$$\widehat{\nabla}_{\text{sym}}^{\text{EP}}(\beta) \triangleq \frac{\widehat{\nabla}^{\text{EP}}(\beta) + \widehat{\nabla}^{\text{EP}}(-\beta)}{2} = -\frac{\partial \mathcal{L}^*}{\partial \theta} + O(\beta^2)$$

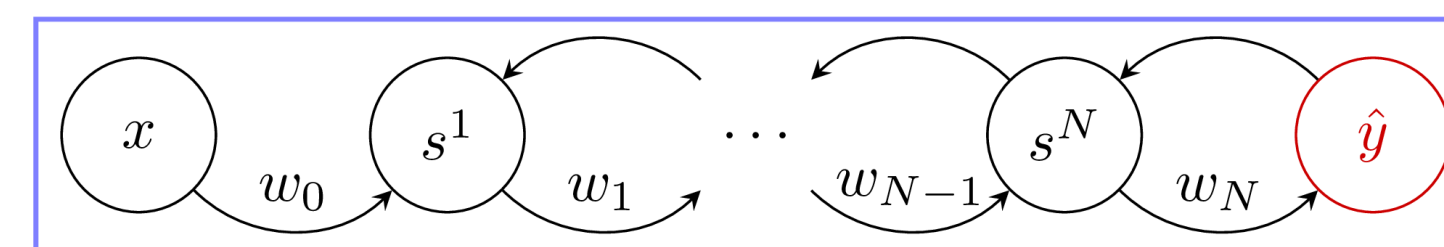
Third phase: $s_*^{-\beta}$



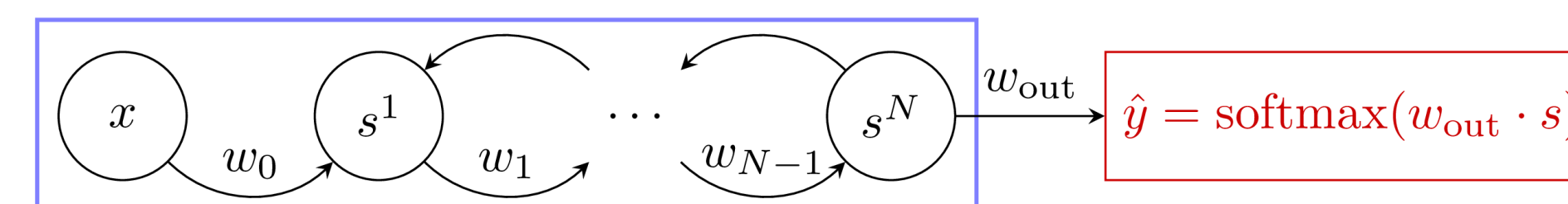
Loss Function	EP Gradient Estimate	EP Error (%)	BPTT Error (%)
Squared Error	2-Phase / $\widehat{\nabla}^{\text{EP}}$	86.7 (5.8)	84.9
	Random Sign	21.6 (20.0)	20.0
	3-Phase / $\widehat{\nabla}_{\text{sym}}^{\text{EP}}$	12.5 (0.2)	7.8

Optimize the cross-entropy loss with EqProp

Output units in the system : difficult to define logits



$$h_t^\beta: \text{hidden layer } \hat{y}_{t+1}^\beta = \frac{\partial \Phi}{\partial \hat{y}}(x, h_t^\beta, \hat{y}_t^\beta, \theta) + \beta (y - \hat{y}_t^\beta)$$



$$s_{t+1}^\beta = \frac{\partial \Phi}{\partial s}(x, s_t^\beta, \theta) + \beta w_{\text{out}}^\top \cdot (y - \hat{y}_t^\beta)$$

Loss Function	EP Gradient Estimate	EP Error (%)	BPTT Error (%)
Cross-Ent.	3-Phase / $\widehat{\nabla}_{\text{sym}}^{\text{EP}}$	11.7 (0.2)	5.0
Cross-Ent. (Dropout)	3-Phase / $\widehat{\nabla}_{\text{sym}}^{\text{EP}}$	11.9 (0.3)	6.5

Distinct forward and backward weights

Generalized framework [3] : $s_{t+1} = F(x, s_t, \theta)$

$$\widehat{\nabla}_{\text{sym}}^{\text{VF}}(\beta) \triangleq \frac{1}{2\beta} \frac{\partial F}{\partial \theta}(x, s_*, \theta)^\top \cdot (s_*^\beta - s_*^{-\beta})$$

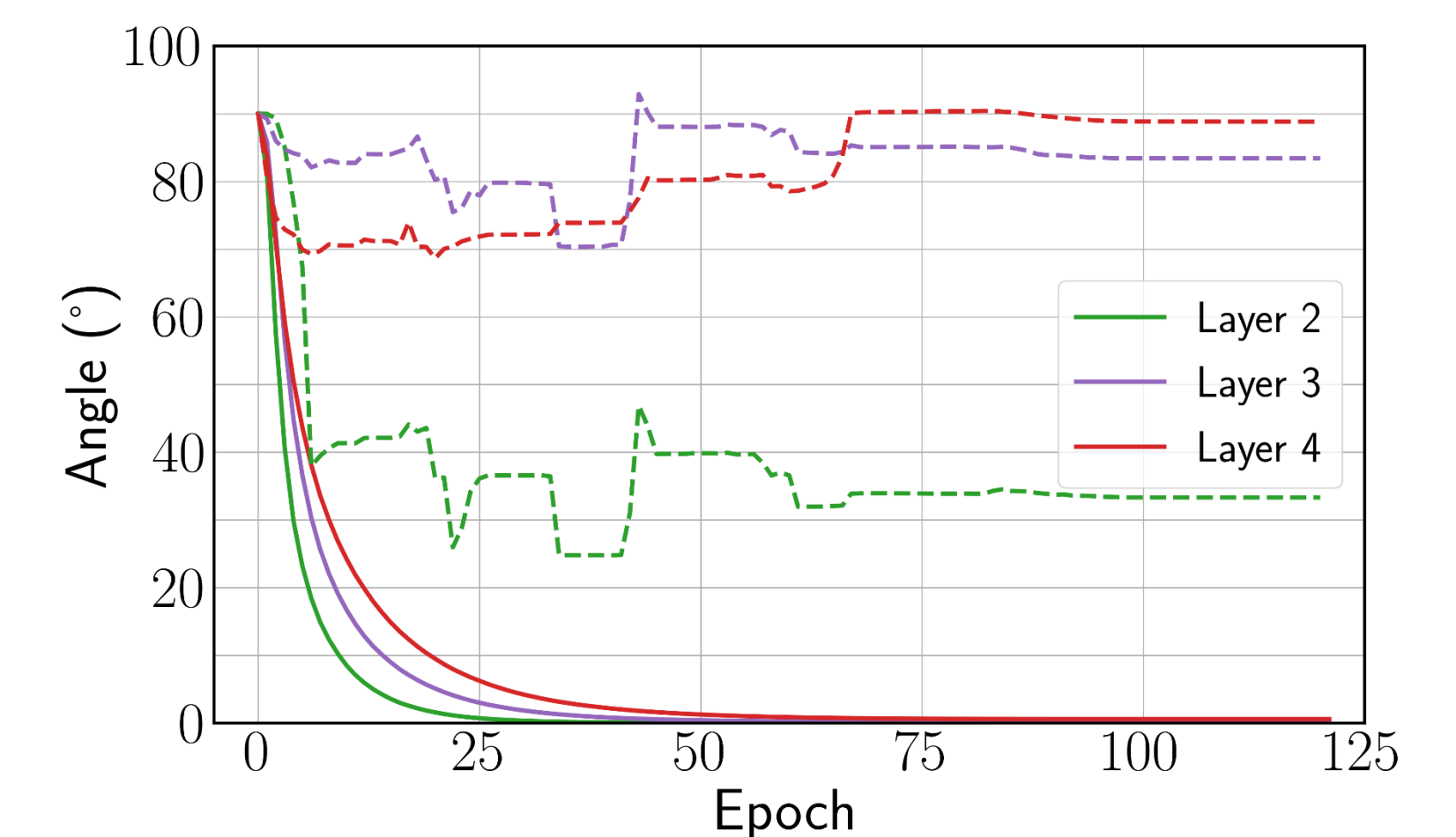
Need for an alignment mechanism, we take inspiration from [4] :

$$\begin{cases} \Delta \theta_f = \eta \left(\widehat{\nabla}_{\text{sym}}^{\text{KP-VF}}(\beta) - \lambda \theta_f \right) \\ \Delta \theta_b = \eta \left(\widehat{\nabla}_{\text{sym}}^{\text{KP-VF}}(\beta) - \lambda \theta_b \right) \end{cases}$$

$$\theta_f(t) - \theta_b(t) = (1 - \eta \lambda)^t (\theta_f(0) - \theta_b(0))$$

Loss Function	EP Gradient Estimate	EP Error (%)	BPTT Error (%)
Cross-Ent.	3-Phase / $\widehat{\nabla}_{\text{sym}}^{\text{VF}}$	75.5 (4.7)	78.0
	3-Phase / $\widehat{\nabla}_{\text{sym}}^{\text{KP-VF}}$	13.2 (0.5)	8.9

We observe that weights do not align (dashed) without the alignment mechanism (solid)



References

- Scellier, B. Bengio, Y. (2017). "Equilibrium propagation: Bridging the gap between energy-based models and backpropagation" In: *Frontiers in computational neuroscience* 11, 24
- Ernoult, M. et al. (2019). "Updates of equilibrium prop match backprop through time in an RNN with static input." In: *Advances in Neural Information Processing Systems* 7081-7091
- Scellier, B. et al. (2018). "Generalization of Equilibrium Propagation to vector fields dynamics." *arXiv:1808.04873*
- Kolen, J. F. and Pollack J. B. (1994). "Backpropagation without weight transport". In: *IEEE ICCN'94*.