

# Fast and deep neuromorphic learning with time-to-first-spike coding

J. Göltz<sup>1</sup>, A. Baumbach<sup>1</sup>, S. Billaudelle<sup>1</sup>, O. Breitwieser<sup>1</sup>, D. Dold<sup>1</sup>, L. Kriener<sup>2</sup>, A. F. Kungl<sup>1</sup>,  
W. Senn<sup>2</sup>, J. Schemmel<sup>1</sup>, K. Meier<sup>1,†</sup>, M. A. Petrovici<sup>2,1</sup>

† Deceased.

<sup>1</sup>Kirchhoff-Institute for Physics, Heidelberg University, 69120 Heidelberg, Germany.

<sup>2</sup>Department of Physiology, University of Bern, 3012 Bern, Switzerland.

## Abstract

*For a biological agent operating under environmental pressure, energy consumption and reaction times are of critical importance. Similarly, engineered systems also strive for short time-to-solution and low energy-to-solution characteristics. At the level of neuronal implementation, this implies achieving the desired results with as few and as early spikes as possible. In the time-to-first-spike coding framework, both of these goals are inherently emerging features of learning. Here, we describe a rigorous derivation of error-backpropagation-based learning for hierarchical networks of leaky integrate-and-fire neurons. We explicitly address two issues that are relevant for both biological plausibility and applicability to neuromorphic substrates by incorporating dynamics with finite time constants and by optimizing the backward pass with respect to substrate variability. This narrows the gap between previous models of first-spike-time learning and biological neuronal dynamics, thereby also enabling fast and energy-efficient inference on analog neuromorphic devices that inherit these dynamics from their biological archetypes, which we demonstrate on two generations of the BrainScaleS analog neuromorphic architecture.*

## 1 Introduction

In recent years, the machine learning landscape has been dominated by deep learning methods. Among the benchmark problems they managed to crack, some were thought to still remain elusive for a long time (LeCun *et al.*, 2015; Krizhevsky *et al.*, 2012; Goodfellow *et al.*, 2014; Silver *et al.*, 2017; Vaswani *et al.*, 2017). It is thus not exaggerated to say that deep learning has reformed our understanding and the future role of “artificial intelligence” (Brooks *et al.*, 2012; Ng, 2016; Hassabis *et al.*, 2017; Sejnowski, 2018; Richards *et al.*, 2019).

However, compared to abstract neural networks used in deep learning, their more biological archetypes – spiking neural networks – still lag behind in performance and scalability (Pfeiffer & Pfeil, 2018). Reasons for this difference in success are numerous; for instance, unlike abstract neu-

rons, even an individual biological neuron represents a complex system, with finite response times, membrane dynamics and spike-based communication (Gerstner, 2001; Izhikevich, 2004), making it more challenging to find reliable coding and computation paradigms (Gerstner, 1998; Maass, 2016; Davies, 2019). Furthermore, one of the major driving forces behind the success of deep learning, the backpropagation of errors algorithm (Rumelhart *et al.*, 1986), remained incompatible with spiking neural networks for a long time (Esser *et al.*, 2015; Schmitt *et al.*, 2017; Tavanaei *et al.*, 2018; Neftci *et al.*, 2019).

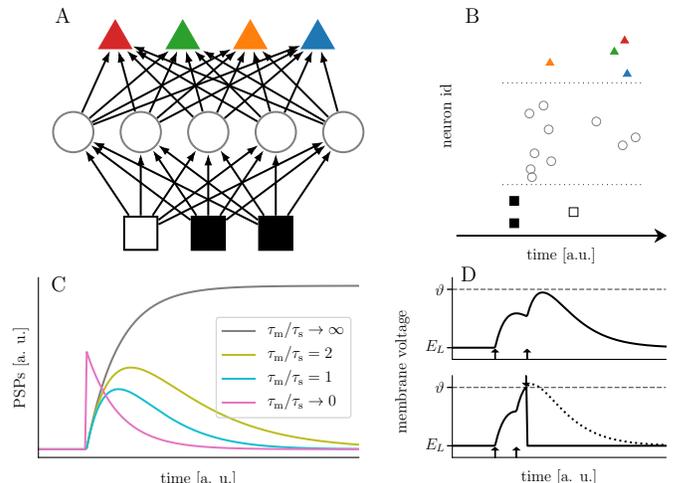
Despite these challenges, spiking neural networks promise to hold some intriguing advantages. The asynchronous nature of spike-based communication allows a coding scheme that utilizes both spatial and temporal dimensions (Gütig & Sompolinsky, 2006), unlike rate-based or spike-count-based approaches (Cao *et al.*, 2015; Diehl *et al.*, 2016; Schmitt *et al.*, 2017; Wu *et al.*, 2019), where the information of spike times is lost due to temporal or population averaging. Due to the inherent parallelism of all biological, as well as many biologically-inspired, neuromorphic systems, this promises fast, sparse and energy-efficient information processing, which might hold the key to novel computing architectures that could one day rival the efficiency of the brain itself (Mead, 1990; Indiveri *et al.*, 2011; Roy *et al.*, 2019). This makes spiking neural networks potentially more powerful than the “conventional”, simple models currently used in machine learning (Maass, 1997), even though this potential still remains mostly unexploited (Pfeiffer & Pfeil, 2018).

Many attempts have been made to reconcile spiking neural networks with their abstract counterparts in terms of functionality, e.g., featuring spike-based stochastic inference models (Petrovici *et al.*, 2013; Neftci *et al.*, 2014; Petrovici *et al.*, 2016; Neftci *et al.*, 2016; Leng *et al.*, 2018; Kungl *et al.*, 2019; Dold *et al.*, 2019; Jordan *et al.*, 2019) and deep models trained on target spike times by shallow learning rules (Kheradpisheh *et al.*, 2018; Illing *et al.*, 2019) or using spike-compatible versions of the error backpropagation algorithm (Bohte *et al.*, 2000; Lee *et al.*, 2016; O’Connor & Welling, 2016; Zenke & Ganguli, 2018; Huh & Sejnowski,

2018; Jin *et al.*, 2018; Tavanaei *et al.*, 2018; Kulkarni & Rajendran, 2018; Wu *et al.*, 2018; Wu *et al.*, 2019; Bellec *et al.*, 2019). A particularly elegant way of utilizing the temporal aspect of exact spike times is the time-to-first-spike (TTFS) coding scheme (Thorpe *et al.*, 2001). Here, a neuron encodes a continuous variable as the time elapsed before its first spike. Such single-spike coding enables fast information processing by inherently encouraging as few and as early spikes as possible, which meets physiological constraints and reaction times observed in humans and animals (Thorpe *et al.*, 1996; Decharms & Merzenich, 1996; Wehr & Laurent, 1996; Johansson & Birznieks, 2004; Gollisch & Meister, 2008; Saal *et al.*, 2016; Portelli *et al.*, 2016). Apart from biological plausibility, such a coding scheme is a natural fit for neuromorphic systems that offer energy-efficient and fast emulation of spiking neural networks (Schemmel *et al.*, 2010; Akopyan *et al.*, 2015; Friedmann *et al.*, 2017; Davies *et al.*, 2018; Mayr *et al.*, 2019; Pei *et al.*, 2019).

For hierarchical TTFS networks, a gradient-descent-based learning rule was proposed in (Mostafa, 2017; Kheradpisheh & Masquelier, 2019), using error backpropagation on a continuous function of output spike times. However, this approach is limited to a neuron model without leak, which is neither biologically plausible, nor compatible with most analog very-large-scale integration (VLSI) neuron dynamics (Thakur *et al.*, 2018). We extend the aforementioned approach to the leaky integrate-and-fire (LIF) model with current-based (CuBa) synapses, which represents an analytically treatable dynamical model of spiking neurons with realistic integration behavior (Rauch *et al.*, 2003; Gerstner & Naud, 2009; Teeter *et al.*, 2018), i.e., with finite membrane ( $\tau_m$ ) and synaptic ( $\tau_s$ ) time constants. For three special cases ( $\tau_m = \tau_s$ ,  $\tau_m = 2\tau_s$  and  $\tau_m \rightarrow \infty$ ), both the times-to-first-spike as well as the gradients of the loss function are analytically calculable.

The closed and exact analytical form of the proposed model, especially for gradients used in weight updates, enables a robust implementation on neuromorphic physical-model systems. We demonstrate such an implementation on the BrainScaleS-2 (Friedmann *et al.*, 2017; Billaudelle *et al.*, 2019) and BrainScaleS-1 (Schemmel *et al.*, 2010; Schmitt *et al.*, 2017; Kungl *et al.*, 2019) accelerated, analog spiking neuromorphic systems. By incorporating information generated on the hardware for updates during training, the algorithm can adapt to the imperfections of the analog circuits. This allows a natural transfer from theory and simulation to a neuromorphic physical model system, demonstrating that the proposed model deals well with various drawbacks of physical-model systems such as fixed-pattern noise and limited parameter precision and control. Such a robustness of coding and learning under imperfections of



**Figure 1: Pattern recognition with time-to-first-spike coding.** (A) Hierarchical feed-forward network structure. Colors encode labels, throughout all figures. (B) Input ( $\square$ ,  $\blacksquare$ ), hidden ( $\circ$ ) and label ( $\blacktriangle$ ) spike times in a raster plot. The TTFS coding amounts to converting black ( $\blacksquare$ )/white ( $\square$ ) pixels to early/late spikes and that the first label neuron to spike determines the inferred class ( $\blacktriangle$ ). (C) Post-synaptic potential (PSP) shapes for different ratios of  $\tau_s$  and  $\tau_m$ . For finite time  $\tau_m$  the membrane ‘forgets’ prior input, making it fundamentally different from the case where  $\tau_m$  is infinite. (D) Illustration of a key challenge posed by finite membrane time constants: small variations of synaptic weights (not shown) or input spike times (upward arrows) can result in an appearing/disappearing output spike and a corresponding discontinuity in the function describing its timing.

the underlying neuronal substrate represents a quintessentially desirable property for every model claiming biological plausibility and for every application geared towards physical computing systems (Prodromakis & Tounazou, 2010; Esser *et al.*, 2015; van De Burgt *et al.*, 2018; Wunderlich *et al.*, 2019; Kungl *et al.*, 2019; Dold *et al.*, 2019; Feldmann *et al.*, 2019).

In the following, we first introduce the CuBa LIF model and the TTFS coding scheme (Section 2.1), before we demonstrate how both inference and training via error backpropagation can be performed analytically with such dynamics (Section 2.2). Finally, the presented model is evaluated both in software simulations (Section 3.1) and neuromorphic emulations (Section 3.2).

## 2 Mathematical results

### 2.1 Preliminaries

**Leaky integrate-and-fire dynamics** The dynamics of an LIF neuron with CuBa synapses are given by

$$C_m \dot{u} = g_l(E_l - u) + \sum_i w_i \sum_{t_i} \theta(t - t_i) \exp\left(-\frac{t - t_i}{\tau_s}\right), \quad (1)$$

with membrane capacitance  $C_m$ , leak conductance  $g_l$ , presynaptic weights  $w_i$  and spike times  $t_i$ , synaptic time constant  $\tau_s$  and  $\theta$  the Heaviside step function. The first sum runs over all presynaptic neurons while the second sum runs over all spikes for each presynaptic neuron. The neuron elicits a spike at time  $T$  when the presynaptic input pushes the membrane potential above a threshold  $\vartheta$ . After spiking, a neuron becomes refractory for a time period  $\tau_{\text{ref}}$ , which is modeled by clamping its membrane potential to a reset value  $\varrho$ :  $u(t') = \varrho$  for  $T \leq t' \leq T + \tau_{\text{ref}}$ . For convenience and without loss of generality, we set the leak potential  $E_l = 0$ . Eqn. (1) can be solved analytically, which yields the sub-threshold dynamics

$$u(t) = \frac{1}{C_m} \frac{\tau_m \tau_s}{\tau_m - \tau_s} \sum_i w_i \kappa(t - t_i), \quad (2)$$

$$\kappa(t) = \theta(t) \left[ \exp\left(-\frac{t}{\tau_m}\right) - \exp\left(-\frac{t}{\tau_s}\right) \right], \quad (3)$$

with membrane time constant  $\tau_m = C_m/g_l$  and the PSP kernel  $\kappa$  given by a difference of exponentials. Here we already assumed the TTFS use case with only one relevant spike for each neuron, for which the second sum in Eqn. (1) reduces to a single term. The choice of  $\tau_m$  ultimately influences the shape of a PSP, starting from a simple exponential ( $\tau_m \ll \tau_s$ ), to a difference of exponentials (with an alpha function for the special case of  $\tau_m = \tau_s$ ) to a graded step function ( $\tau_m \gg \tau_s$ ) (Fig. 1C).

The first two cases are markedly different from the last one, which is also known as either the non-leaky integrate-and-fire (nLIF) or simply integrate-and-fire (IF) model and was used in previous work (Mostafa, 2017). In the nLIF model, input to the membrane is never forgotten, as opposed to the LIF model, where the PSP reaches a peak after finite time and subsequently decays back to its baseline. In other words, presynaptic spikes in the LIF model have a purely local effect in time, unlike in the nLIF model, where only the onset of a PSP is localized in time, but the postsynaptic effect remains forever, or until the postsynaptic neuron spikes. A finite  $\tau_m$  thus assigns much more importance to the time differences between input spikes and introduces discontinuities in the neuronal output that make an analytical treatment more difficult (Fig. 1D).

**Time-to-first-spike coding** Our spike-based neural code follows an idea first proposed in (Mostafa, 2017). Unlike coding in artificial neural networks (ANNs) and different from rate-based codes in spiking neural networks (SNNs), this scheme explicitly uses the timing of individual spikes for encoding information. In time-to-first-spike (TTFS) coding, the presence of a feature is reflected by

the timing of a neuron’s first spike, with earlier spikes representing a more strongly manifested feature. This has the effect that important information inherently propagates faster through the network, with potentially only few spikes needed for the network to process an input. Consequently, this scheme enables a more efficient processing of inputs, both in terms of time-to-solution and energy-to-solution (assuming the latter depends on the total number of spikes and the time required for the network to solve, e.g., an input classification problem).

## 2.2 Learning rules

In order to formulate the optimization of first-spike times  $T$  as a gradient-descent problem, we need to derive closed-form expressions for these  $T$ . This is equivalent to finding the time of the first threshold crossing by solving  $u(T) = \vartheta$  for  $T$ . Even though a general closed-form solution does not exist, we show analytical solutions for three specific cases: (i)  $\tau_m \rightarrow \infty$ , (ii)  $\tau_m = \tau_s$  and (iii)  $\tau_m = 2\tau_s$ :

$$\frac{T}{\tau_s} = \ln \left[ \frac{a_1}{a_\infty - \vartheta C_m / \tau_s} \right], \quad \tau_m = \infty \text{ (nLIF)}; \quad (4)$$

$$\frac{T}{\tau_s} = \frac{b}{a_1} - \mathcal{W} \left[ \underbrace{-\frac{g_l \vartheta}{a_1} \exp\left(\frac{b}{a_1}\right)}_{=: z} \right], \quad \tau_m = \tau_s \text{ (LIF)}; \quad (5)$$

$$\frac{T}{\tau_s} = 2 \ln \left[ \frac{2a_1}{a_2 + \sqrt{a_2^2 - 4a_1 g_l \vartheta}} \right], \quad \tau_m = 2\tau_s \text{ (LIF)}; \quad (6)$$

where  $\mathcal{W}$  is the Lambert W function and using

$$a_n := \sum_{i \in C} w_i e^{t_i / n \tau_s}, \quad (7)$$

$$b := \sum_{i \in C} w_i \frac{t_i}{\tau_s} e^{t_i / \tau_s}, \quad (8)$$

as shorthand for sums over the set of causal presynaptic spikes  $C = \{i \mid t_i < T\}$ . All three equations are differentiable with respect to synaptic weights and presynaptic spike times. For a detailed derivation of these results, we refer to Appendix A.

The implicit assumption of having only the first spike emitted by every neuron be relevant for downstream processing can effectively be ensured by using a long enough refractory period. Since the only information-carrying signal that is not reset upon firing is the synaptic memory, which is forgotten on the time scale of  $\tau_s$ , we found that, in practice, setting  $\tau_{\text{ref}} = \tau_s$  leads to most neurons eliciting only one spike before the classification of a given input pattern.

The case  $\tau_m \rightarrow \infty$  has already been discussed in Mostafa (2017) and was reproduced here for completeness and comparison. Due to the symmetry in  $\tau_m$  and  $\tau_s$  of the PSP (Eqn. A7), the  $\tau_m = 2\tau_s$  case describes the  $\tau_m = \frac{1}{2}\tau_s$  case as well. Using Eqns. (4) to (6), we can treat the TTFS network much like an ANN, where instead of rates, spike times are propagated. In a layered feed-forward network, Eqns. (4) to (6) can be used iteratively, i.e., one can calculate the spike times of the first layer, use these to calculate the spike times of the second layer, etc., until the label neurons are reached.

While we found both rules for finite  $\tau_m$  to work well in practice, we focus on  $\tau_m = \tau_s$  in the following, as  $\tau_m = 2\tau_s$  and  $\tau_m = \frac{1}{2}\tau_s$  can be treated analogously. Equations for all cases are derived in Appendix A.

**Exact error backpropagation with spikes** As depicted in Fig. 1A, we consider feed-forward networks of CuBa LIF neurons. The input uses the same coding scheme as all other neurons, with input neurons spiking earlier for darker pixels. In particular, no external time reference is required: the network effectively processes contrast information and is essentially agnostic with respect to specific absolute input spike times. The output of the network is defined by the identity of the label neuron that spikes first (Fig. 1B).

We denote by  $t_k^{(l)}$  the output spike time of the  $k$ th neuron in the  $l$ th layer, e.g., for a network with  $N$  layers,  $t_k^{(N)}$  is the spike time of the  $k$ th neuron in the label layer. The weight projecting to the  $k$ th neuron of layer  $l$  from the  $i$ th neuron of layer  $l-1$  is denoted by  $w_{ki}^{(l)}$ .

To apply a variant of the error backpropagation algorithm (Rumelhart *et al.*, 1986), we choose a loss function that is differentiable with respect to synaptic weights and spike times. During learning, the objective is to maximize the temporal difference between the correct and all other label spikes while minimizing the time-to-correct-solution. The following loss function fulfills the above requirements:

$$L[\mathbf{t}^{(N)}, p] = -\log \left[ \frac{\exp(-t_p^{(N)}/\xi\tau_s)}{\sum_k \exp(-t_k^{(N)}/\xi\tau_s)} \right] + \alpha \left[ \exp\left(\frac{t_p^{(N)}}{\beta\tau_s}\right) - 1 \right], \quad (9)$$

where  $\mathbf{t}^{(N)}$  denotes the vector of label spike times  $t_k^{(N)}$ ,  $p$  the index of the correct label and  $\xi$ ,  $\alpha$  and  $\beta$  represent scaling parameters. Because the softmax-scaled spike times can be viewed as assigning a probability to the different labels, the first term represents a cross-entropy of this distribution

relative to the true label distribution (which is 1 for the correct label and 0 otherwise). Reducing this term therefore increases the temporal difference between the output spike of the correct label neuron and all other label neurons. Notably, it only depends on the spike time difference and is invariant under absolute time shifts, making it independent of artificial outside clocking. The second term is a regularizer that favors solutions where the correct label neuron spikes as early as possible.

Weights are updated such that they minimize the loss  $L[\mathbf{t}^{(N)}, p]$ . For weights projecting into the label layer, updates are calculated via

$$\Delta w_{ki}^{(N)} \propto -\frac{\partial L[\mathbf{t}^{(N)}, p]}{\partial w_{ki}^{(N)}} = -\frac{\partial t_k^{(N)}}{\partial w_{ki}^{(N)}} \frac{\partial L[\mathbf{t}^{(N)}, p]}{\partial t_k^{(N)}}, \quad (10)$$

where the second term can be obtained straightforwardly from Eqn. (9): The first term depends on the PSP shape; the corresponding differentiation of Eqn. (5) results in

$$\frac{\partial t_k^{(l)}}{\partial w_{ki}^{(l)}}(\mathbf{w}^{(l)}, \mathbf{t}^{(l-1)}) = -\frac{1}{a_1} \exp\left[\frac{t_i^{(l-1)}}{\tau_s}\right] \times \left[ \left(\frac{t_i^{(l-1)}}{\tau_s} - \frac{b}{a_1}\right)(1 - z\mathcal{W}'(z)) + z\mathcal{W}'(z) \right], \quad (11)$$

for an arbitrary layer  $l$ , where  $z$  and  $a_1$  are given in Eqns. (5) and (7). Using a relation for the derivative of  $\mathcal{W}$ , the equation can be simplified and made to depend on the output spike time  $\mathbf{t}^{(l)}$ :

$$\frac{\partial t_k^{(l)}}{\partial w_{ki}^{(l)}}(\mathbf{w}, \mathbf{t}^{(l-1)}, \mathbf{t}^{(l)}) = -\frac{1}{a_1} \exp\left[\frac{t_i^{(l-1)}}{\tau_s}\right] \frac{1}{W(z) + 1} \frac{t_k^{(l)} - t_i^{(l-1)}}{\tau_s}. \quad (12)$$

Using this additional information optimizes learning in scenarios where the inferred output spike and the true output spike differ (Section 3.1).

The weight updates of deeper layers can be calculated iteratively by application of the chain rule:

$$\Delta w_{ki}^{(l)} \propto -\frac{\partial L[\mathbf{t}^{(N)}, p]}{\partial w_{ki}^{(l)}} = -\frac{\partial t_k^{(l)}}{\partial w_{ki}^{(l)}} \delta_k^{(l)}, \quad (13)$$

where the second term is a propagated error that can be calculated recursively with a sum over the neurons in layer  $(l+1)$ :

$$\delta_k^{(l)} := \frac{\partial L[\mathbf{t}^{(N)}, p]}{\partial t_k^{(l)}} = \sum_j \frac{\partial t_j^{(l+1)}}{\partial t_k^{(l)}} \delta_j^{(l+1)}. \quad (14)$$

The latter derivative amounts to, once the output spike time is reinserted as above,

$$\begin{aligned} & \frac{\partial t_k^{(l)}}{\partial t_i^{(l-1)}}(\mathbf{w}, \mathbf{t}^{(l-1)}, \mathbf{t}^{(l)}) \\ &= -\frac{1}{a_1} \exp\left[\frac{t_i^{(l-1)}}{\tau_s}\right] \frac{1}{W(z)+1} \frac{w_{ki}^{(l)} t_k^{(l)} - t_i^{(l-1)} - \tau_s}{\tau_s}. \end{aligned} \quad (15)$$

The learning rule can be rewritten in layer-wise form to resemble the standard error backpropagation algorithm for abstract neurons (see Appendix B for the standard backpropagation equation):

$$\delta^{(N)} = \frac{\partial L}{\partial \mathbf{t}^{(N)}}, \quad (16)$$

$$\delta^{(l-1)} = \boldsymbol{\rho}^{(l-1)} \odot \left( \widetilde{\mathbf{W}}^{(l)T} \delta^{(l)} \right), \quad (17)$$

$$\Delta \mathbf{w}^{(l)} = -\eta \left( \delta^{(l)} \boldsymbol{\rho}^{(l-1)T} \right) \odot \widehat{\mathbf{B}}^{(l)}, \quad (18)$$

where  $\odot$  is the element-wise product, the  $T$ -superscript denotes the transpose of a matrix and  $\delta^{(l-1)}$  is a vector containing the backpropagated errors of layer  $(l-1)$ . The individual elements of the tensors above are given by

$$\rho_i^{(l)} = -\frac{1}{a_1} \exp\left[\frac{t_i^{(l)}}{\tau_s}\right] \frac{1}{W(z)+1} \frac{1}{\tau_s}, \quad (19)$$

$$\widehat{B}_{ki}^{(l)} = t_k^{(l)} - t_i^{(l-1)}, \quad (20)$$

$$\widetilde{W}_{ki}^{(l)} = w_{ki}^{(l)} \frac{t_k^{(l)} - t_i^{(l-1)} - \tau_s}{\tau_s}. \quad (21)$$

In this form, it becomes apparent that for training, only the label layer error and the neuron spike times are required, which can either be calculated using Eqn. (5) or by simulating (or emulating) the LIF dynamics (Eqn. 1).

As mentioned above, the treatment of the other two special cases is analogous to the above. Thus, for CuBa LIF neurons with finite time constants  $\tau_m = \tau_s$ ,  $\tau_m = 2\tau_s$  and  $\tau_m = \frac{1}{2}\tau_s$ , both the forward pathway (spike times) and backward pathway (backpropagation) can be calculated analytically using a loss that is differentiable with respect to both synaptic weights and neuronal spike times.

## 3 Demonstrating learning on various spiking substrates

### 3.1 Simulation results

**Classification task** We demonstrate the above framework in a pattern classification task (Fig. 2A), with the

spiking network (Fig. 1A) simulated in NEST (Gewaltig & Diesmann, 2007). To assist learning, mini batch training was used and the weight updates  $\Delta w$  calculated as described in Section 2.2 were L1-normalized layer-wise to be smaller than 10. Furthermore, for layers with more than 35% of silent neurons averaged over the minibatches, all afferent weights were increased by a fixed amount in order to have sufficient activity. In case of multiple layers where this applied, only the first layer with insufficient spikes was boosted.

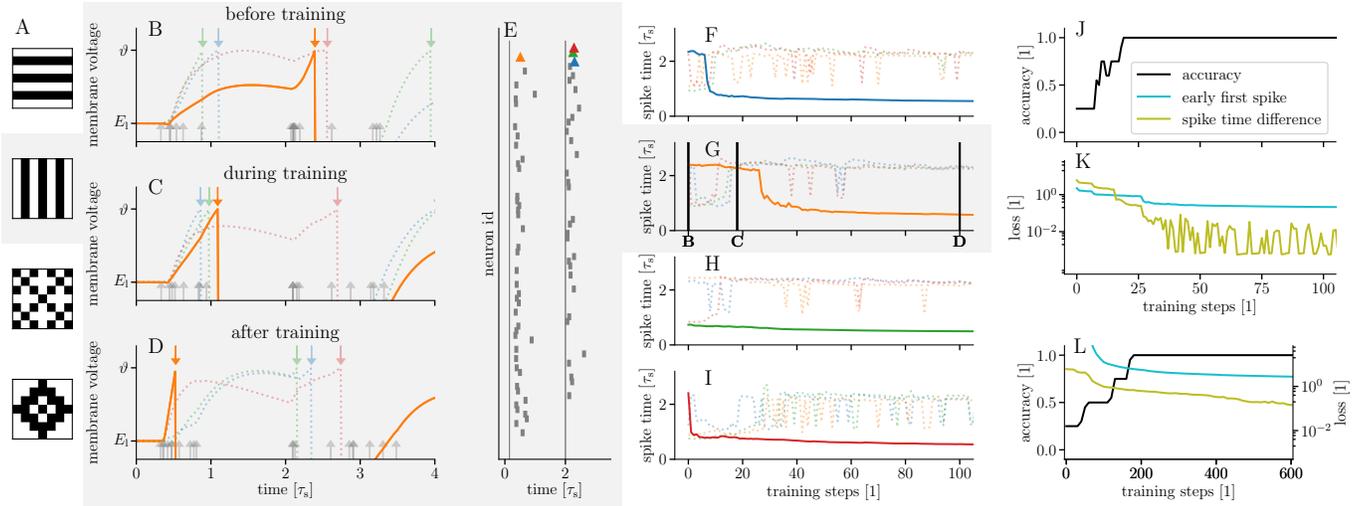
Fig. 2B-K shows results from training a network with 49 visible, 80 hidden and 4 label neurons on this data set. While not used during training, the temporal evolution of the membrane potentials helps illustrate the learning process. Fig. 2B-D shows voltages in the label layer for one class (orange) before, during and after training, illustrating how the trained weights make the correct neuron spike earliest by a large margin (see also Fig. 2E).

The spike times including the input spikes (vertical lines) and the ones in the hidden layer are shown in a raster plot in Fig. 2E. Due to the finite membrane and synaptic time constants, output spikes can only happen within a finite time window after their inputs. In the particular training scenario described in Fig. 2, this renders output spikes happening immediately before the late input spikes extremely unlikely. This effective gap can also be seen in the evolution of the spike times during training, where a small change in synaptic weights can bring a spike from the early into the late group and vice-versa, leading to sudden jumps in both spike timing and loss (Fig. 2F-I,K; see also Fig. 1D).

The evolution of the label layer spike times for all four classes is shown in Fig. 2F-I, including the steps at which the voltage plots were recorded. The spike times for the different classes together decide both the accuracy (the proportion of correct classifications) and the loss (Eqn. 9). The evolution of both during training is shown in Fig. 2J,K.

It is important to note that we have chosen a simple dataset in order to make it amenable to emulation on our neuromorphic systems (Section 3.2). In particular, it is linearly separable and would thus not require backpropagation for perfect classification. Therefore, to demonstrate that error backpropagation is working as intended, we performed an additional simulation with frozen weights between the hidden and label layer, training only the ones between input and hidden layer. As expected and shown in Fig. 2L, training was successful in this setup, but took considerably longer for the same initialization and hyperparameters as used in Fig. 2B-K.

As mentioned above, our loss function consists of two parts (Eqn. 9), the first relating to TTFS coding and the second representing a regularizer that pushes correct neurons



**Figure 2: Pattern recognition with time-to-first-spike coding for a simple data set.** (A) Input pattern set consisting of four classes. (B-D) Voltage dynamics in the label layer (colored traces) and output spikes (downward arrows) induced by spikes from the hidden layer (upward arrows) for one pattern (plots of data from the same pattern are marked by gray background) at different stages of training. As intended, the learning rule has the effect of decreasing the spike time of the correct label neuron ( $\blacktriangledown$ ) while increasing the spike time of the other neurons. The colors of the traces correspond to the four different label neurons (Fig. 1A,B), with the correct one shown in orange. All times are given in units of the synaptic time constant  $\tau_s$ . (E) Raster plot after training for the same sample, including spikes in the hidden layer (gray marks) and early and late input times (vertical lines). The classification (first label spike, orange) happens prior to a significant fraction of the hidden neuron spikes. (F-I) Evolution of label neuron spike times during training for all four classes, with colors marking the different label neurons as above, and incorrect labels being lighter. The correct neuron’s spike time decreases while all others are pushed back, producing a distinct gap. In (G) the snapshots from B, C, and D are marked. (J) Evolution of accuracy during training. (K) Loss (green is the first term of Eqn. (9), blue the  $\alpha$ -weighted second term) only reaches small values once the accuracy (I) is already at 100%. (L) To show that back-propagation is working we trained only the weights from input to hidden neurons, keeping those from the hidden to label neurons fixed.

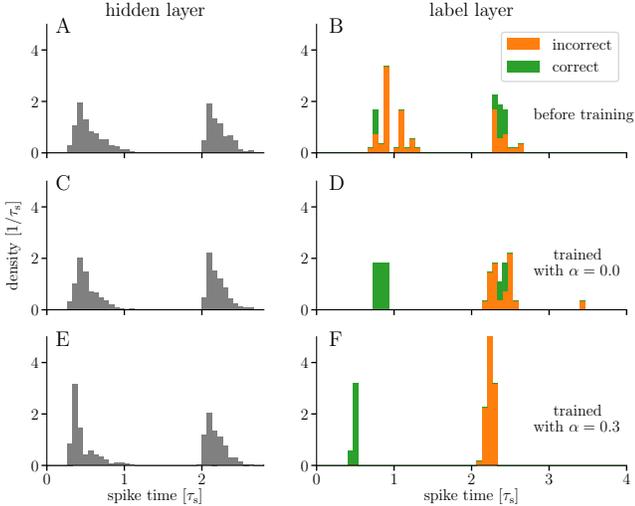
to spike early, and stabilises the training. Figure 3 shows the effect of regularization: it shifts correct label spikes to earlier times, which in turn causes the afferent active hidden neurons to spike earlier as well.

**Inserting substrate-specific information into the backward pass** As noted in the introduction, TTFS coding is a natural fit for neuromorphic hardware due to its emphasis on speed (early  $\tau_s$ ), especially for accelerated devices which can profit additionally from intrinsic speed-up (Section 3.2). However, this speed comes at the price of reduced control over certain neuron and synapse parameters. This implies in particular that the ratio of the membrane and synaptic time constants is different from the ideal values of  $\frac{\tau_m}{\tau_s} = 1$  or  $\frac{\tau_m}{\tau_s} = 2$  used in derivations. It is therefore crucial that the learning rule still works under such parameter variability in order for it to be applicable to such neuromorphic substrates. The question at hand touches upon whether our learning rules also work for other parameter values than the specific ones for which they were derived, and if so, how well. We study the robustness of learning for different time constant ratios by sweeping a range of membrane time constants using the NEST (Gewaltig & Diesmann, 2007) simulator for the forward pass and the different learning rules for the backward pass.

One detail is important in this context. Eqns. (4) to (6) for the spike times depend only on neuron parameters, presynaptic spike times and weights, thus the derivatives needed for the weight update initially depend on those ‘natural’ variables as well (Eqns. A15 and A22). With some manipulations, the equation for the actual output spike time can be inserted (Eqns. A17 and A24), producing a version of the learning rule that profits from more information of the forward pass and is thus significantly more stable. The two versions are identical only in case the forward and backward pass have exactly the same ratios. The effect of this disparity is shown in Fig. 4. For both update rules, including information about the true spiking activity significantly improves learning over a wide range of  $\tau_m/\tau_s$  ratios.

### 3.2 Fast neuromorphic classification

In this framework, classification speed is a function of the network depth and the time constants  $\tau_m$  and  $\tau_s$ . Assuming typical biological timescales, most input patterns in the above scenario are classified within several ms. By leveraging the speedup of neuromorphic systems such as BrainScaleS (Schmitt *et al.*, 2017; Billaudelle *et al.*, 2019), with intrinsic acceleration factors of  $10^3$ - $10^4$ , the same computation can be achieved within  $\mu$ s.

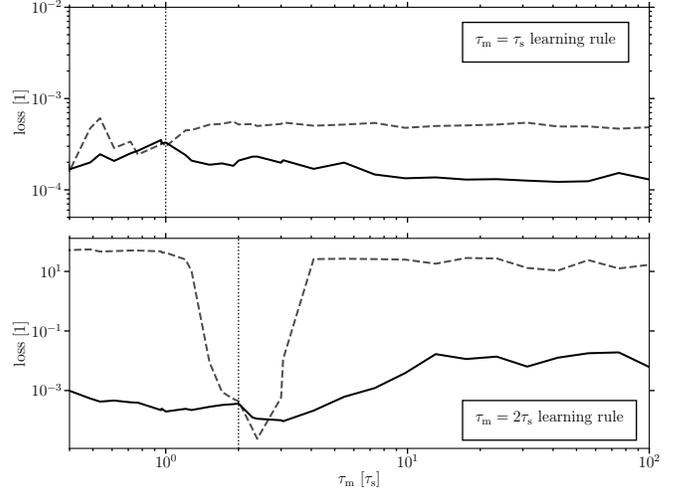


**Figure 3: Spike time distributions** over all four classes before (A,B) and after training without ( $\alpha = 0$ , C, D) and with ( $\alpha = 0.3$ , E, F) regularization. Left column (gray): hidden layer; right column: label layer, with correct spikes marked in green and false ones in orange. Here, we used noised inputs, with five examples per class, i.e. 20 patterns in total. The separation of the distribution into two distinct modes is a direct consequence of the black/white input and its encoding (see also Fig. 2). While the network trains to 100% accuracy for both values of  $\alpha$  (data for  $\alpha = 0$  not shown), a nonzero  $\alpha$  leads to significantly earlier spikes in the label layer, as well as to the correct label neurons never spiking during the second volley.

However, the speed advantages of such analog systems compared to software simulations come at the cost of reduced control, and training needs to cope with phenomena such as spike time jitter and neuron parameter variability. In particular, this implies  $\tau_m \neq \tau_s$ , so the derived learning rule is only an approximation of true gradient descent in these systems, as discussed above.

### 3.2.1 Learning with TTFS on BrainScaleS-2

We ported the network architecture and training scheme to the BrainScaleS-2 system, a mixed-signal accelerated neuromorphic platform. The application-specific integrated circuit (ASIC) is built around an analog neuromorphic core which emulates the dynamics of neurons and synapses. All state variables, such as membrane potentials and synaptic currents, are physically represented in their respective circuits and evolve continuously in time. Considering the natural time constants of such integrated analog circuits, this emulation takes place at 1000-fold accelerated time scales compared to the biological nervous system. One BrainScaleS-2 core features 512 AdEx neurons, which can be freely configured; these circuits can be restricted to LIF dynamics as required by our training framework (Aamir *et al.*, 2018b; Aamir *et al.*, 2018a). Both the membrane and synaptic time constants were calibrated to 5  $\mu$ s.

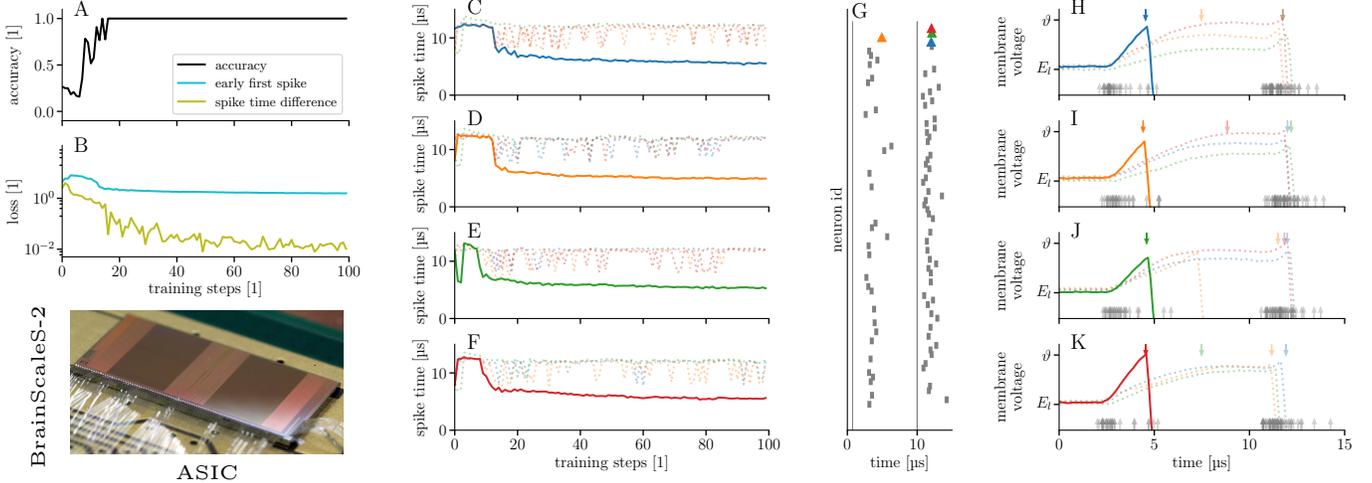


**Figure 4: Fitness of learning rules for varying  $\tau_m$ .** We simulated the forward pass for different  $\tau_m/\tau_s$  ratios by varying the  $\tau_m$  of neurons in the network. Here, we show the median cross-entropy over the last 300 training steps after a total of 3000 training steps, averaged over 30 different random initializations. This allows us to compare the efficacy of learning in two scenarios for the backward pass. Solid lines: weight updates depending only on neuron parameters, presynaptic spike times and weights (Eqn. 11). Dashed lines: weight updates including the true output spike times (Eqns. 15 and 12). Dotted lines: ideal  $\tau_m/\tau_s$  ratios. Note how including output spike times significantly improves the training for both learning rules across a wide range of parameters.

Each neuron circuit integrates stimuli from a column of 256 current-based synapses (Friedmann *et al.*, 2017). Each synapse holds a 6 bit weight value; its sign is shared with all other synapses located on the same row in the synapse matrix. The presented training scheme, however, allows weights to continuously transition between excitation and inhibition. We therefore allocated pairs of synapse rows to convey the activity of single presynaptic partners, one configured for excitation, the other one for inhibition.

Synapses receive their inputs from an event routing module allowing to connect neurons within a chip as well as to inject stimuli from external sources. Events emitted by the neuron circuits are annotated with a time stamp and then sent off-chip. The neuromorphic ASIC is accompanied by a field-programmable gate array (FPGA) to handle the communication with the host computer. It also provides mechanisms for low-latency experiment control including the timed release of spike trains into the neuromorphic core. The FPGA is furthermore used to record events and digitized membrane traces originating from the ASIC.

We used an in-the-loop-training approach, where emulation runs on the neuromorphic substrate were interleaved with host-based weight update calculations (Schmitt *et al.*, 2017). For the emulation of the forward pass, the data set was broken down into mini-batches, converted into input spike trains and then injected into the neuromorphic sys-



**Figure 5: Neuromorphic pattern recognition with time-to-first-spike coding on BrainScaleS-2.** (A) Accuracy and (B) energy (green: cross-entropy, blue:  $\alpha$ -weighted regularizer in Eqn. (9)) during training. (C-F) Evolution of label neuron spike times, displayed separately for the four classes. For each pattern, the spike time of the neuron representing this class is shown as a solid line in full color. (G) Raster plot of hidden neurons (gray) and the four label neurons (colored) after training, shown for a stimulus representing the second pattern ( $\blacktriangle$ ). (H-K) Membrane voltage traces of the label neurons, for the four classes respectively. These analog membrane traces were digitized on the neuromorphic substrate after 100 training steps.

tem via the FPGA. The latter was also used to record the spikes emitted by the hidden and label layers. Weight updates were – based on these output spike trains – calculated on the host computer and then written back to the synapse memory. This backward pass shared its implementation with the previously described simulation framework.

We were able to successfully and reliably train the network emulated on BrainScaleS-2 on the discussed data set (Fig. 5). The system quickly learned to fully discriminate between the presented patterns, with clear separation between label spike times. Learning performance in terms of convergence speed is difficult to compare because the hyperparameters are not easily transferable, but appears similar to the numerical simulations of the same network. After training, due to the interplay of the system’s intrinsic acceleration and the nature of the learning algorithm itself, each pattern was classified in less than 5  $\mu$ s.

### 3.2.2 Learning with TTFS on BrainScaleS-1

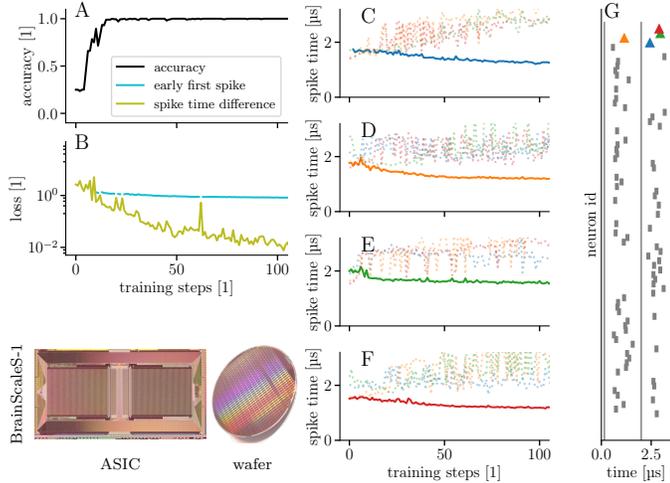
To demonstrate the amenability of our approach to different neuromorphic substrates, we also tested it on the BrainScaleS-1 system (Schemmel *et al.*, 2010). This version of BrainScaleS has a very similar architecture to BrainScaleS-2, but its component chips are interconnected through postprocessing on their common wafer (wafer-scale integration). More importantly for our coding scheme and learning rules, its circuits emulate conductance-based (CoBa) instead of CuBa neurons. Furthermore, due to the different fabrication technology and design choices (in par-

ticular, the floating-gate parameter memory, see Srowig *et al.*, 2007; Schemmel *et al.*, 2010; Koke, 2017), the parameter variability and spike time jitter are significantly higher than on BrainScaleS-2 (Schmitt *et al.*, 2017).

The training procedure was analogous to the one used on BrainScaleS-2. To accommodate the CoBa synapse dynamics, we introduced global weight scale factors that modeled the distance between reversal and leak potentials and the total conductance, which were multiplied to the synaptic weights to achieve a CuBa approximation for which our learning rules apply. Despite this approximation and the considerable substrate variability (compare, e.g., Fig. 6C-F with Fig. 5C-F), our framework was able to compensate well, almost matching the performance achieved on BrainScaleS-2 (Fig. 6).

## 4 Discussion

In this manuscript, we proposed a model of deep time-to-first-spike learning that builds on a principled view of neurosynaptic dynamics with finite time constants and comes with exact learning rules for optimizing first-spike times; an early version of this work was presented in Göltz (2019). In this quintessentially spike-based learning framework, only single spike times are required for calculating the weight updates, thus reducing the memory (bandwidth and capacity) requirements of synaptic updates in comparison to, e.g., rate coding approaches (see, e.g., Schmitt *et al.*, 2017, for an example of deep but rate-based learning that was also



**Figure 6: Training a spiking network on the wafer-scale BrainScaleS-1 system.** Accuracy (A) and loss (B) during training of the four pattern data set (Fig. 2A). Green: cross-entropy, blue:  $\alpha$ -weighted regularizer in Eqn. (9). (C-F) Evolution of the spike times in the label layer for the four different patterns. In each, the neuron coding the correct class is shown with a solid line and in full color. (G) Raster plot after training for the second pattern ( $\blacktriangle$ ).

applied to a BrainScaleS system).

Our work builds on earlier results by Mostafa (2017), which we extended to accommodate leaky integrate-and-fire neurons, thereby including more biologically plausible and neuromorphic-hardware-compatible neuro-synaptic dynamics. Additionally, we introduced a regularizing loss term that favors early classification, thereby significantly improving the time-to-solution of the network. To account for substrate variability, we further incorporated output spike times directly into the backward pass, which extends the applicability of our derived learning rules to a wide range of parameters, thus allowing us to demonstrate the framework on two different neuromorphic platforms (two generations of the BrainScaleS architecture) that exhibit varying degrees of parameter noise in their analog components. Unlike other approaches (Mostafa, 2017; Comsa *et al.*, 2019; Kheradpisheh & Masquelier, 2019) we do not use any kind of clocking or bias spikes, thereby being independent of any absolute time reference or global clock signal.

The complexity of the learned dataset was mostly limited by the size of the used substrate, but we expect the framework to scale to significantly more challenging problems, as suggested by the FPGA-based experiments in Mostafa *et al.* (2017). After learning, the network needed less than one spike per neuron to produce a correct classification on all used substrates. With these few spikes, we achieved a time-to-solution of less than one synaptic time constant. Since the dynamical timescales directly affect the duration of the network emulation between synaptic updates, this inher-

ently leads to a significant reduction of the total training time. Taking into consideration relaxation times between patterns, our setup was able to handle a concatenated pattern throughput of at least 20 kHz, independently of emulated network size. These results promise an efficient exploitation of such accelerated neuromorphic substrates for high-throughput inference on spiking input data.

## Acknowledgment

The authors wish to thank Sebastian Schmitt for BrainScaleS-1 support as well as Jakob Jordan and Nico Görtler for valuable discussions.

**Funding** The authors gratefully acknowledge funding from the European Union under grant agreements 604102, 720270, 785907 (HBP) and the Manfred Stärk Foundation.

**Simulation** Simulations were performed on the bw-ForCluster NEMO, supported by the state of Baden-Württemberg through bwHPC and the German Research Foundation (DFG) through grant no INST 39/963-1 FUGG.

## Appendix

### A Derivation of main results

In this section we derive the equations in the main manuscript, starting with the learning rule for  $\tau_m \rightarrow \infty$ , Eqn. (4), then  $\tau_m = \tau_s$ , Eqn. (5) and finally  $\tau_m = 2\tau_s$ , Eqn. (6).

For each, a solution for the spike time  $T$ , defined by

$$u(T) = \vartheta, \quad (\text{A1})$$

given LIF dynamics

$$u(t) = \frac{1}{C_m} \frac{\tau_m \tau_s}{\tau_m - \tau_s} \sum_{\text{spikes } t_i} w_i \kappa(t - t_i), \quad (\text{A2})$$

$$\kappa(t) = \theta(t) \left[ \exp\left(-\frac{t}{\tau_m}\right) - \exp\left(-\frac{t}{\tau_s}\right) \right], \quad (\text{A3})$$

has to be found. For convenience, we use the following definitions

$$a_n := \sum_{i \in C} w_i e^{t_i/n\tau_s}, \quad (\text{A4})$$

$$b := \sum_{i \in C} w_i \frac{t_i}{\tau_s} e^{t_i/\tau_s}, \quad (\text{A5})$$

with summation over the set of causal presynaptic spikes  $C = \{i \mid t_i < T\}$ .

## A.1 nLIF learning rule for $\tau_m \rightarrow \infty$

With this choice of  $\tau_m$ , the first term in Eqn. (A2) becomes 1 and we recover the nLIF case discussed in (Mostafa, 2017). Given the existence of an output spike, in Eqn. (A1) the spike time  $T$  appears only in one place and simple reordering yields

$$\frac{T}{\tau_s} = \ln \left[ \frac{a_1}{a_\infty - \vartheta C_m / \tau_s} \right], \quad (\text{A6})$$

where we used Eqn. (A4) for  $n = 1$  and  $n = \infty$ , the latter being the sum over the weights.

## A.2 Learning rule for $\tau_m = \tau_s$

**Spike time** According to l'Hôpital's rule, in the limit  $\tau_m \rightarrow \tau_s$  Eqn. (A2) becomes a sum over  $\alpha$ -functions of the form

$$u(t) = \frac{1}{C_m} \sum_i w_i \theta(t - t_i) \cdot (t - t_i) \exp \left( -\frac{t - t_i}{\tau_s} \right). \quad (\text{A7})$$

Using these voltage dynamics for the equation of the spike time Eqn. (A1), together with the definition Eqn. (A5) and  $\tau_m = C_m / g_l$ , we get the equation

$$0 = g_l \vartheta \exp \left[ \frac{T}{\tau_s} \right] + \underbrace{b - a_1 \frac{T}{\tau_s}}_{=: y}. \quad (\text{A8})$$

The variable  $y$  is introduced to bring the equation into the form

$$h e^h = z \quad (\text{A9})$$

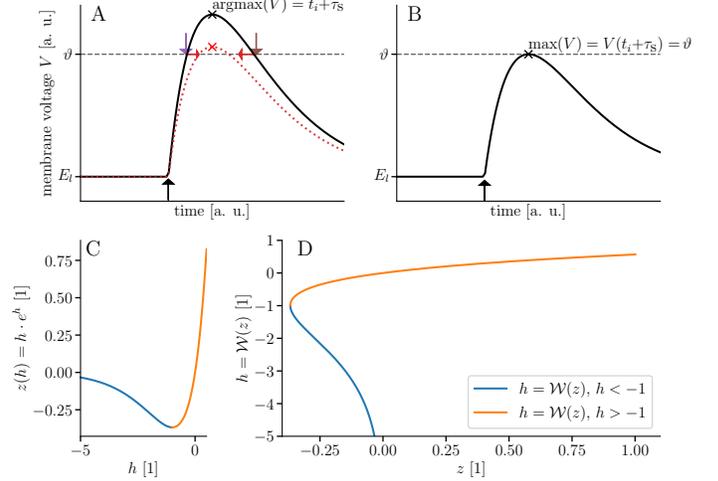
which can be solved with the differentiable Lambert W function  $h = \mathcal{W}(z)$ . The goal is now to bring Eqn. (A8) into this form, this is achieved by reformulation in terms of  $y$

$$0 = g_l \vartheta \exp \left( \frac{b}{a_1} \right) \exp \left( -\frac{y}{a_1} \right) + y \quad (\text{A10})$$

$$\underbrace{\frac{y}{a_1}}_{=: h} \exp \left( \frac{y}{a_1} \right) = - \underbrace{\frac{g_l \vartheta}{a_1} \exp \left( \frac{b}{a_1} \right)}_{=: z}. \quad (\text{A11})$$

With the definition of the Lambert W function the spike time can be written as

$$\frac{T}{\tau_s} = \frac{b}{a_1} - \mathcal{W} \left[ -\frac{g_l \vartheta}{a_1} \exp \left( \frac{b}{a_1} \right) \right]. \quad (\text{A12})$$



**Figure A1:** (A) Membrane dynamics for one strong input spike at  $t_i$  (upward arrow) with two threshold crossings due to leak pullback (earlier violet, later brown). The change induced by a reduction of the input weight is shown in red. (B) Edge case without crossing and exactly one time where  $V(t) = \vartheta$ . (C) Defining relation for the Lambert W function  $\mathcal{W}$ , evidently not an injective map. (D) Distinguishing between  $h \leq -1$  allows to define the inverse function of (C), the Lambert W function  $\mathcal{W}$ .

**Branch choice** Given that a spike happens, there will be two threshold crossings: One from below at the actual spike time, and one from above when the voltage decays back to the leak potential (Fig. A1A,B). Correspondingly, the Lambert W function (Fig. A1C,D) has two real branches (in addition to infinite imaginary ones), and we need to choose the branch that returns the earlier solution. In case the voltage is only tangent to the threshold at its maximum, the Lambert W function only has one solution.

For choosing the branch in the other cases we need to look at  $h$  from the definition, i.e.

$$h = \frac{y}{a_1} = \frac{b}{a_1} - \frac{T}{\tau_s}. \quad (\text{A13})$$

In a setting with only one strong enough input spike, the summations in  $a_n$  and  $b$  reduce to yield  $h = (t_i - T) / \tau_s$ . Because the maximum of the PSP for  $\tau_m = \tau_s$  occurs at  $t_i + \tau_s$ , we know that the spike must occur at  $T \leq t_i + \tau_s$  and therefore

$$-1 \leq \frac{t_i - T}{\tau_s} = h. \quad (\text{A14})$$

This corresponds to the branch cut of the Lambert W function meaning we must choose the branch with  $h \geq -1$ . For a general setting, if we know a spike exists, we expect  $a_n$  and  $b$  to be positive. In order to get the earlier threshold crossing, we need the branch that returns the larger  $\mathcal{W}$  (Fig. A1D), that is where  $\mathcal{W} = h > -1$ .

**Derivatives** The derivatives for  $t_i$  in the causal set  $i \in C$  come down to

$$\frac{\partial T}{\partial t_i}(\mathbf{w}, \mathbf{t}) \quad (\text{A15})$$

$$= -\frac{1}{a_1} \exp\left[\frac{t_i}{\tau_s}\right] \frac{w_i}{\tau_s} \left[1 + \left(\frac{t_i}{\tau_s} - \frac{b}{a_1}\right) (1 - z\mathcal{W}'(z))\right],$$

$$\frac{\partial T}{\partial w_i}(\mathbf{w}, \mathbf{t}) \quad (\text{A16})$$

$$= -\frac{1}{a_1} \exp\left[\frac{t_i}{\tau_s}\right] \left[z\mathcal{W}'(z) + \left(\frac{t_i}{\tau_s} - \frac{b}{a_1}\right) (1 - z\mathcal{W}'(z))\right].$$

A crucial step is to reinsert the definition of the spike time where it is possible (cf. Section 3.1). For this we need the derivative of the Lambert W function  $z\mathcal{W}'(z) = \frac{\mathcal{W}(z)}{\mathcal{W}(z)+1}$  that follows from differentiating its definition Eqn. (A9) with  $h = \mathcal{W}(z)$  w.r.t.  $z$ . With this derivative one can calculate the derivative of Eqn. (A12) with respect to incoming weights and times as functions of presynaptic weights, input spike times and output spike time:

$$\frac{\partial T}{\partial t_i}(\mathbf{w}, \mathbf{t}, T) = -\frac{1}{a_1} \frac{1}{\mathcal{W}(z)+1} \exp\left[\frac{t_i}{\tau_s}\right] \frac{w_i}{\tau_s} \frac{T - t_i - \tau_s}{\tau_s}, \quad (\text{A17})$$

$$\frac{\partial T}{\partial w_i}(\mathbf{w}, \mathbf{t}, T) = -\frac{1}{a_1} \frac{1}{\mathcal{W}(z)+1} \exp\left[\frac{t_i}{\tau_s}\right] \frac{T - t_i}{\tau_s}. \quad (\text{A18})$$

### A.3 Learning rule for $\tau_m = 2\tau_s$

**Spike time** Inserting the voltage (Eqn. A2) into the spike time (Eqn. A1) yields

$$g_1\vartheta = e^{-T/\tau_m} \sum_{i \in C} w_i e^{\frac{t_i}{\tau_m}} - e^{-T/\tau_s} \sum_{i \in C} w_i e^{\frac{t_i}{\tau_s}}. \quad (\text{A19})$$

Reordering and rewriting this in terms of  $a_1$ ,  $a_2$ , and  $\tau_s$  (with  $\tau_m = 2\tau_s$ ) we get

$$0 = -a_1 \left(e^{-T/2\tau_s}\right)^2 + a_2 e^{-T/2\tau_s} - g_1\vartheta. \quad (\text{A20})$$

This is written such that its quadratic nature becomes apparent, making it possible to solve for  $\exp(-T/2\tau_s)$  and thus

$$\frac{T}{\tau_s} = 2 \ln \left[ \frac{2a_1}{a_2 + \sqrt{a_2^2 - 4a_1 g_1\vartheta}} \right]. \quad (\text{A21})$$

**Branch choice** The quadratic equation has two solutions that correspond to the voltage crossing at spike time and relaxation towards the leak later; again, we want the earlier of the two solutions. It follows from the monotonicity of the logarithm that the earlier time is the one with the larger denominator. Due to an output spike requiring an excess of

recent positively weighted input spikes,  $a_n$  are positive, and the + solution is the correct one.

**Derivatives** Using the definition  $x = \sqrt{a_2^2 - 4a_1 g_1\vartheta}$  for brevity, the derivatives of Eqn. (A21) are

$$\frac{\partial T}{\partial w_i}(\mathbf{w}, \mathbf{t}) \quad (\text{A22})$$

$$= 2\tau_s \left[ \frac{1}{a_1} + \frac{2g_1\vartheta}{(a_2+x)x} \right] \exp\left[\frac{t_i}{\tau_s}\right] - \frac{2\tau_s}{x} \exp\left[\frac{t_i}{2\tau_s}\right],$$

$$\frac{\partial T}{\partial t_i}(\mathbf{w}, \mathbf{t}) \quad (\text{A23})$$

$$= 2 \left[ \frac{1}{a_1} + \frac{2g_1\vartheta}{(a_2+x)x} \right] \exp\left[\frac{t_i}{\tau_s}\right] - \frac{1}{x} \exp\left[\frac{t_i}{2\tau_s}\right].$$

Again, inserting the output spike time yields

$$\frac{\partial T}{\partial w_i}(\mathbf{w}, \mathbf{t}, T) \quad (\text{A24})$$

$$= \frac{2\tau_s}{a_1} \left[ 1 + \frac{g_1\vartheta}{x} \exp\left[\frac{T}{2\tau_s}\right] \right] \exp\left[\frac{t_i}{\tau_s}\right] - \frac{2\tau_s}{x} \exp\left[\frac{t_i}{2\tau_s}\right],$$

$$\frac{\partial T}{\partial t_i}(\mathbf{w}, \mathbf{t}, T) \quad (\text{A25})$$

$$= \frac{2}{a_1} \left[ 1 + \frac{g_1\vartheta}{x} \exp\left[\frac{T}{2\tau_s}\right] \right] \exp\left[\frac{t_i}{\tau_s}\right] - \frac{1}{x} \exp\left[\frac{t_i}{2\tau_s}\right].$$

## B Standard error backpropagation

The standard error backpropagation formula for artificial (rate-based) neural networks (Rumelhart *et al.*, 1986) with rates  $\mathbf{a}$  is given by

$$\delta^{(N)} = \frac{\partial L}{\partial \mathbf{a}^{(N)}} \odot \mathbf{a}'^{(N)}, \quad (\text{A26})$$

$$\delta^{(l-1)} = \mathbf{a}'^{(l-1)} \odot \left( \mathbf{W}^{(l)T} \delta^{(l)} \right), \quad (\text{A27})$$

$$\Delta \mathbf{W}^{(l)} = -\eta \delta^{(l)} \mathbf{a}^{(l-1)T}. \quad (\text{A28})$$

Traditionally, in artificial neural networks, the last layer is a linear classifier, but here, to highlight the resemblance to rate-based neurons, we define the loss function on the activation of the neurons in the last layer  $L = \frac{1}{2} \|\mathbf{y} - \mathbf{a}^N\|^2$ , where  $\mathbf{y}$  is the target label in one-hot coding.

## References

1. Aamir, S. A. *et al.* A Mixed-Signal Structured AdEx Neuron for Accelerated Neuromorphic Cores. *IEEE Transactions on Biomedical Circuits and Systems* **12**, 1027–1037 (2018).

2. Aamir, S. A. *et al.* An Accelerated LIF Neuronal Network Array for a Large-Scale Mixed-Signal Neuromorphic Architecture. *IEEE Transactions on Circuits and Systems I: Regular Papers* **65**, 4299–4312 (2018).
3. Akopyan, F. *et al.* Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **34**, 1537–1557 (2015).
4. Bellec, G. *et al.* Biologically inspired alternatives to backpropagation through time for learning in recurrent neural nets. *CoRR* **abs/1901.09049**. arXiv: 1901 . 09049. <<http://arxiv.org/abs/1901.09049>> (2019).
5. Billaudelle, S. *et al.* Versatile emulation of spiking neural networks on an accelerated neuromorphic substrate. in preparation (2019).
6. Bohte, S. M., Kok, J. N. & La Poutré, J. A. *Spike-Prop: backpropagation for networks of spiking neurons*. in *ESANN* (2000), 419–424.
7. Brooks, R., Hassabis, D., Bray, D. & Shashua, A. Is the brain a good model for machine intelligence? *Nature* **482**, 462 (2012).
8. Cao, Y., Chen, Y. & Khosla, D. Spiking deep convolutional neural networks for energy-efficient object recognition. *International Journal of Computer Vision* **113**, 54–66 (2015).
9. Comsa, I. M. *et al.* Temporal coding in spiking neural networks with alpha synaptic function. *arXiv preprint arXiv:1907.13223* (2019).
10. Davies, M. Benchmarks for progress in neuromorphic computing. *Nature Machine Intelligence* **1**, 386–388 (2019).
11. Davies, M. *et al.* Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro* **38**, 82–99 (2018).
12. Decharms, R. C. & Merzenich, M. M. Primary cortical representation of sounds by the coordination of action-potential timing. *Nature* **381**, 610 (1996).
13. Diehl, P. U., Zarella, G., Cassidy, A., Pedroni, B. U. & Neftci, E. *Conversion of artificial recurrent neural networks to spiking neural networks for low-power neuromorphic hardware* in *2016 IEEE International Conference on Rebooting Computing (ICRC)* (2016), 1–8.
14. Dold, D. *et al.* Stochasticity from function—Why the Bayesian brain may need no noise. *Neural Networks* **119**, 200–213 (2019).
15. Esser, S. K., Appuswamy, R., Merolla, P., Arthur, J. V. & Modha, D. S. *Backpropagation for energy-efficient neuromorphic computing* in *Advances in Neural Information Processing Systems* (2015), 1117–1125.
16. Feldmann, J., Youngblood, N., Wright, C., Bhaskaran, H & Pernice, W. All-optical spiking neurosynaptic networks with self-learning capabilities. *Nature* **569**, 208 (2019).
17. Friedmann, S. *et al.* Demonstrating Hybrid Learning in a Flexible Neuromorphic Hardware System. *IEEE Transactions on Biomedical Circuits and Systems* **11**, 128–142. ISSN: 1932-4545 (2017).
18. Gerstner, W. *Spiking neurons* tech. rep. (MIT-press, 1998).
19. Gerstner, W. in *Plausible neural networks for biological modelling* 23–48 (Springer, 2001).
20. Gerstner, W. & Naud, R. How good are neuron models? *Science* **326**, 379–380 (2009).
21. Gewaltig, M.-O. & Diesmann, M. NEST (NEural Simulation Tool). *Scholarpedia* **2**, 1430 (2007).
22. Gollisch, T. & Meister, M. Rapid neural coding in the retina with relative spike latencies. *science* **319**, 1108–1111 (2008).
23. Göltz, J. *Training Deep Networks with Time-to-First-Spike Coding on the BrainScaleS Wafer-Scale System* Master’s thesis (Universität Heidelberg, Apr. 2019). <<http://www.kip.uni-heidelberg.de/Veroeffentlichungen/details.php?id=3909>>.
24. Goodfellow, I. *et al.* *Generative adversarial nets* in *Advances in neural information processing systems* (2014), 2672–2680.
25. Gütig, R. & Sompolinsky, H. The tempotron: a neuron that learns spike timing-based decisions. *Nature neuroscience* **9**, 420 (2006).
26. Hassabis, D., Kumaran, D., Summerfield, C. & Botvinick, M. Neuroscience-inspired artificial intelligence. *Neuron* **95**, 245–258 (2017).
27. Huh, D. & Sejnowski, T. J. in *Advances in Neural Information Processing Systems 31* (eds Bengio, S. *et al.*) 1433–1443 (Curran Associates, Inc., 2018). <<http://papers.nips.cc/paper/7417-gradient-descent-for-spiking-neural-networks.pdf>>.
28. Illing, B., Gerstner, W. & Brea, J. Biologically plausible deep learning—but how far can we go with shallow networks? *Neural Networks* (2019).
29. Indiveri, G. *et al.* Neuromorphic silicon neuron circuits. *Frontiers in neuroscience* **5**, 73 (2011).

30. Izhikevich, E. M. Which model to use for cortical spiking neurons? *IEEE transactions on neural networks* **15**, 1063–1070 (2004).
31. Jin, Y., Zhang, W. & Li, P. *Hybrid macro/micro level backpropagation for training deep spiking neural networks* in *Advances in Neural Information Processing Systems* (2018), 7005–7015.
32. Johansson, R. S. & Birznieks, I. First spikes in ensembles of human tactile afferents code complex spatial fingertip events. *Nature neuroscience* **7**, 170 (2004).
33. Jordan, J. *et al.* Deterministic networks for probabilistic computing. *Scientific Reports* **9**, 1–17 (2019).
34. Kheradpisheh, S. R., Ganjtabesh, M., Thorpe, S. J. & Masquelier, T. STDP-based spiking deep convolutional neural networks for object recognition. *Neural Networks* **99**, 56–67 (2018).
35. Kheradpisheh, S. R. & Masquelier, T. S4NN: temporal backpropagation for spiking neural networks with one spike per neuron. *arXiv preprint arXiv:1910.09495* (2019).
36. Koke, C. *Device Variability in Synapses of Neuromorphic Circuits* PhD thesis (2017). doi:10.11588/heidok.00022742.
37. Krizhevsky, A., Sutskever, I. & Hinton, G. E. *Imagenet classification with deep convolutional neural networks* in *Advances in neural information processing systems* (2012), 1097–1105.
38. Kulkarni, S. R. & Rajendran, B. Spiking neural networks for handwritten digit recognition—Supervised learning and network optimization. *Neural Networks* **103**, 118–127 (2018).
39. Kungl, A. F. *et al.* Accelerated physical emulation of Bayesian inference in spiking neural networks. *Frontiers in Neuroscience* **13**, 1201 (2019).
40. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *nature* **521**, 436–444 (2015).
41. Lee, J. H., Delbruck, T. & Pfeiffer, M. Training deep spiking neural networks using backpropagation. *Frontiers in neuroscience* **10**, 508 (2016).
42. Leng, L. *et al.* Spiking neurons with short-term synaptic plasticity form superior generative networks. *Scientific reports* **8**, 10651 (2018).
43. Maass, W. Networks of spiking neurons: the third generation of neural network models. *Neural networks* **10**, 1659–1671 (1997).
44. Maass, W. Searching for principles of brain computation. *Current Opinion in Behavioral Sciences* **11**, 81–92 (2016).
45. Mayr, C., Hoepfner, S. & Furber, S. SpiNNaker 2: A 10 Million Core Processor System for Brain Simulation and Machine Learning. *arXiv preprint arXiv:1911.02385* (2019).
46. Mead, C. Neuromorphic electronic systems. *Proceedings of the IEEE* **78**, 1629–1636 (1990).
47. Mostafa, H., Pedroni, B. U., Sheik, S. & Cauwenberghs, G. *Fast classification using sparsely active spiking networks* in *2017 IEEE International Symposium on Circuits and Systems (ISCAS)* (May 2017), 1–4. doi:10.1109/ISCAS.2017.8050527.
48. Mostafa, H. Supervised learning based on temporal coding in spiking neural networks. *IEEE transactions on neural networks and learning systems* **29**, 3227–3235 (2017).
49. Neftci, E., Das, S., Pedroni, B., Kreuz-Delgado, K. & Cauwenberghs, G. Event-driven contrastive divergence for spiking neuromorphic systems. *Frontiers in neuroscience* **7**, 272 (2014).
50. Neftci, E. O., Mostafa, H. & Zenke, F. Surrogate gradient learning in spiking neural networks. *arXiv preprint arXiv:1901.09948* (2019).
51. Neftci, E. O., Pedroni, B. U., Joshi, S., Al-Shedivat, M. & Cauwenberghs, G. Stochastic synapses enable efficient brain-inspired learning machines. *Frontiers in neuroscience* **10**, 241 (2016).
52. Ng, A. What artificial intelligence can and can't do right now. *Harvard Business Review* **9** (2016).
53. O'Connor, P. & Welling, M. Deep spiking networks. *arXiv preprint arXiv:1602.08323* (2016).
54. Pei, J. *et al.* Towards artificial general intelligence with hybrid Tianjic chip architecture. *Nature* **572**, 106–111 (2019).
55. Petrovici, M. A., Bill, J., Bytschok, I., Schemmel, J. & Meier, K. Stochastic inference with deterministic spiking neurons. *arXiv preprint arXiv:1311.3211* (2013).
56. Petrovici, M. A., Bill, J., Bytschok, I., Schemmel, J. & Meier, K. Stochastic inference with spiking neurons in the high-conductance state. *Physical Review E* **94**, 042312 (2016).
57. Pfeiffer, M. & Pfeil, T. Deep learning with spiking neurons: opportunities and challenges. *Frontiers in neuroscience* **12** (2018).
58. Portelli, G. *et al.* Rank order coding: a retinal information decoding strategy revealed by large-scale multielectrode array retinal recordings. *Eneuro* **3** (2016).

59. Prodrumakis, T. & Toumazou, C. *A review on memristive devices and applications in 2010 17th IEEE International Conference on Electronics, Circuits and Systems* (2010), 934–937.
60. Rauch, A., La Camera, G., Luscher, H.-R., Senn, W. & Fusi, S. Neocortical pyramidal cells respond as integrate-and-fire neurons to in vivo-like input currents. *Journal of neurophysiology* **90**, 1598–1612 (2003).
61. Richards, B. A. *et al.* A deep learning framework for neuroscience. *Nature neuroscience* **22**, 1761–1770 (2019).
62. Roy, K., Jaiswal, A. & Panda, P. Towards spike-based machine intelligence with neuromorphic computing. *Nature* **575**, 607–617 (2019).
63. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by back-propagating errors. *Nature* **533–536**, 533–536 (1986).
64. Saal, H. P., Wang, X. & Bensmaia, S. J. Importance of spike timing in touch: an analogy with hearing? *Current Opinion in Neurobiology* **40**, 142–149 (2016).
65. Schemmel, J. *et al.* *A wafer-scale neuromorphic hardware system for large-scale neural modeling in Proceedings of 2010 IEEE International Symposium on Circuits and Systems* (2010), 1947–1950.
66. Schmitt, S. *et al.* *Neuromorphic hardware in the loop: Training a deep spiking network on the brainscales wafer-scale system in 2017 International Joint Conference on Neural Networks (IJCNN)* (2017), 2227–2234.
67. Sejnowski, T. J. *The deep learning revolution* (MIT Press, 2018).
68. Silver, D. *et al.* Mastering the game of go without human knowledge. *Nature* **550**, 354 (2017).
69. Srowig, A. *et al.* Analog floating gate memory in a 0.18  $\mu\text{m}$  single-poly CMOS process. *Internal FACETS documentation* (2007).
70. Tavanaei, A., Ghodrati, M., Kheradpisheh, S. R., Masquelier, T. & Maida, A. Deep learning in spiking neural networks. *Neural Networks* (2018).
71. Teeter, C. *et al.* Generalized leaky integrate-and-fire models classify multiple neuron types. *Nature communications* **9**, 709 (2018).
72. Thakur, C. S. T. *et al.* Large-scale neuromorphic spiking array processors: A quest to mimic the brain. *Frontiers in neuroscience* **12**, 891 (2018).
73. Thorpe, S., Delorme, A. & Van Rullen, R. Spike-based strategies for rapid processing. *Neural networks* **14**, 715–725 (2001).
74. Thorpe, S., Fize, D. & Marlot, C. Speed of processing in the human visual system. *nature* **381**, 520 (1996).
75. Van De Burgt, Y., Melianas, A., Keene, S. T., Malliaras, G. & Salleo, A. Organic electronics for neuromorphic computing. *Nature Electronics* **1**, 386–397 (2018).
76. Vaswani, A. *et al.* *Attention is all you need in Advances in neural information processing systems* (2017), 5998–6008.
77. Wehr, M. & Laurent, G. Odour encoding by temporal sequences of firing in oscillating neural assemblies. *Nature* **384**, 162 (1996).
78. Wu, J. *et al.* Deep Spiking Neural Network with Spike Count based Learning Rule. *CoRR* **abs/1902.05705**. arXiv: 1902.05705. <<http://arxiv.org/abs/1902.05705>> (2019).
79. Wu, Y., Deng, L., Li, G., Zhu, J. & Shi, L. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in neuroscience* **12** (2018).
80. Wunderlich, T. *et al.* Demonstrating advantages of neuromorphic computation: a pilot study. *Frontiers in Neuroscience* **13**, 260 (2019).
81. Zenke, F. & Ganguli, S. Superspike: Supervised learning in multilayer spiking neural networks. *Neural computation* **30**, 1514–1541 (2018).