

A solution to the learning dilemma for recurrent networks of spiking neurons

Guillaume Bellec^{1,◦}, Franz Scherr^{1,◦}, Anand Subramoney¹, Elias Hajek¹,
Darjan Salaj¹, Robert Legenstein¹ & Wolfgang Maass^{1,*}

¹*Institute of Theoretical Computer Science, Graz University of Technology,
Inffeldgasse 16b, Graz, Austria*

[◦] *Equal contributions.*

^{*} *To whom correspondence should be addressed; E-mail: maass@igi.tugraz.at.*

Abstract

Recurrently connected networks of spiking neurons underlie the astounding information processing capabilities of the brain. But in spite of extensive research, it has remained open how they can learn through synaptic plasticity to carry out complex network computations. We argue that two pieces of this puzzle were provided by experimental data from neuroscience. A new mathematical insight tells us how these pieces need to be combined to enable biologically plausible online network learning through gradient descent, in particular deep reinforcement learning. This new learning method – called *e-prop* – approaches the performance of *BPTT* (backpropagation through time), the best known method for training recurrent neural networks in machine learning. In addition, it suggests a method for powerful on-chip learning in novel energy-efficient spike-based hardware for AI.

Introduction

Networks of neurons in the brain differ in at least two essential aspects from deep neural networks in machine learning: They are recurrently connected, forming a giant number of loops, and they communicate via asynchronously emitted stereotypical electrical pulses, called spikes, rather than bits or numbers that are produced in a synchronized manner by each layer of a feedforward deep network. Models that capture primary information processing capabilities of spiking neurons in the brain are well known, and we consider the arguably most prominent one: leaky integrate-and-fire (LIF) neurons, where spikes that arrive from other neurons through synaptic connections are multiplied with the corresponding synaptic weight, and are linearly integrated by a leaky membrane potential. The neuron fires – i.e., emits a spike – when the membrane potential reaches a firing threshold.

32 But it is an open problem how recurrent networks of spiking neurons (RSNNs) can learn,
33 i.e., how their synaptic weights can be modified by local rules for synaptic plasticity so that
34 the computational performance of the network improves. In deep learning this problem
35 is solved for feedforward networks through gradient descent for a loss function E that
36 measures imperfections of current network performance [1]. Gradients of E are propagated
37 backwards through all layers of the feedforward network to each synapse through a process
38 called backpropagation. Recurrently connected networks can compute more efficiently
39 because each neuron can participate several times in a network computation, and they are
40 able to solve tasks that require integration of information over time or a non-trivial timing
41 of network outputs according to task demands. But since each synaptic weight can affect
42 the network computation at several time points during a recurrent network computation,
43 its impact on the loss function (see Fig. 1a) is more indirect, and learning through gradient
44 descent becomes substantially more difficult. This learning problem is aggravated if there
45 are slowly changing hidden variables in the neuron model, such as neurons with spike-
46 frequency adaptation (SFA). Neurons with SFA are quite common in the neocortex [2],
47 and it turns out that their inclusion in the RSNN significantly increases the computational
48 power of the network [3]. In fact, RSNNs trained through gradient descent acquire then
49 similar computing capabilities as networks of LSTM (Long Short-Term Memory) units,
50 the state of the art for recurrent neural networks in machine learning. Because of this
51 functional relation to LSTM networks these RSNN models are referred to as LSNNs [3].

52 In machine learning one trains recurrent neural networks by unrolling the network into
53 a virtual feedforward network [1], see Fig. 1b, and applying the backpropagation algorithm
54 to that (Fig. 1c). This learning method for recurrent neural networks is called backpropa-
55 gation through time (*BPTT*) since it requires propagation of gradients backwards in time
56 with regard to the network computation.

57 With a careful choice of the pseudo-derivative for handling the discontinuous dynamics
58 of spiking neurons one can apply *BPTT* also to RSNNs, and RSNNs were able to learn in
59 this way for the first time to solve really demanding computational tasks (see [3], [4] for
60 preceding results). But the dilemma is that *BPTT* requires storing the intermediate states
61 of all neurons during a network computation, and merging these in a subsequent offline
62 process with gradients that are computed backwards in time (see Fig. 1c, Movie S1 and
63 Movie S2). This makes it very unlikely that *BPTT* is used by the brain [5].

64 We present a solution to this dilemma in the form of a biologically plausible method for
65 online network learning through gradient descent: *e-prop* (Fig. 1d, see Movie S3). *E-prop*
66 is motivated by two streams of experimental data from neuroscience:

- 67 i) Neurons in the brain maintain traces of preceding activity on the molecular level, for
68 example in the form of calcium ions or activated CaMKII enzymes [6]. In particular,
69 they maintain a fading memory of events where the presynaptic neuron fired before
70 the postsynaptic neuron, which is known to induce synaptic plasticity if followed by
71 a top-down learning signal [7, 8, 9]. Such traces are often referred to as eligibility
72 traces.
- 73 ii) In the brain there exists an abundance of top-down signals such as dopamine, acetyl-

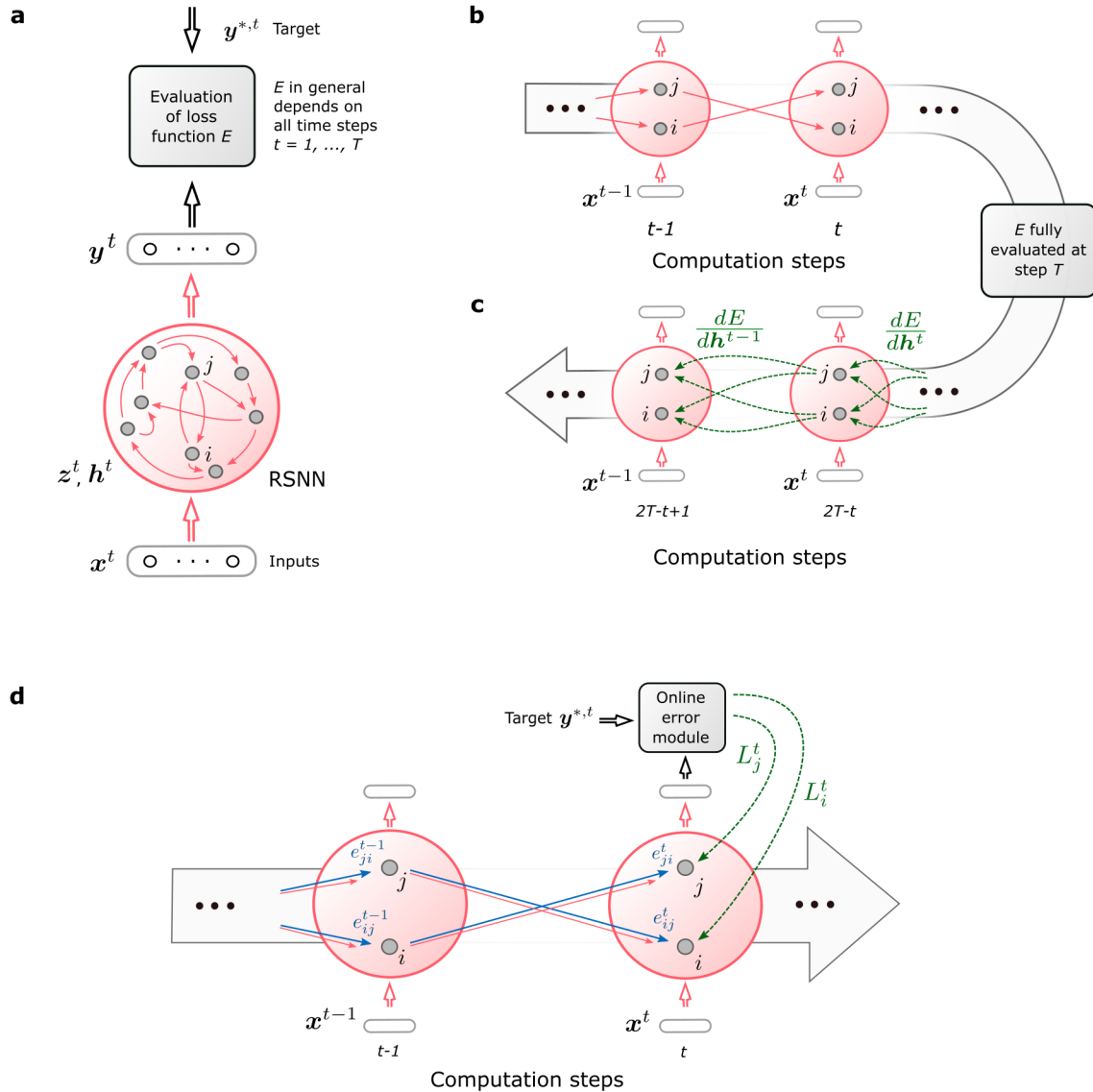


Figure 1: **Schemes for RSNNs, BPTT, and e-prop.** **a)** RSNN with network inputs x , neuron spikes z , and output targets y^* , for each time step t of the RSNN computation. Output neurons y provide a low-pass filter of a weighted sum of network spikes z . **b)** BPTT computes gradients in the unrolled version of the network. It has a new copy of the neurons of the RSNN for each time step t . A synaptic connection from neuron i to neuron j of the RSNN is replaced by an array of feedforward connections, one for each time step t , that goes from the copy of neuron i in the layer for time step t to a copy of neuron j in the layer for time step $t + 1$. All synapses in this array have the same weight: the weight of this synaptic connection in the RSNN. **c)** Loss gradients of BPTT are propagated backwards in time and retrograde across synapses in an offline manner, long after the forward computation has passed a layer. **d)** Online learning dynamics of e-prop. Feedforward computation of eligibility traces is indicated in blue. These are combined with online learning signals according to equation (1).

74 choline, and neural firing [10] related to the event-related negativity (ERN), that
75 inform local populations of neurons about behavioral results. Furthermore dopamine
76 signals [11, 12] have been found to be specific for different target populations of neu-
77 rons, rather than being global. We refer in our learning model to such top-down
78 signals as learning signals.

79 A re-analysis of the mathematical basis of gradient descent in recurrent neural networks
80 tells us how local eligibility traces and top-down learning signals should be optimally com-
81 bined to enable network learning through gradient descent – without requiring backpro-
82 gation of signals through time. The resulting new learning method, *e-prop*, learns slower
83 than *BPTT*, but tends to approximate the performance of *BPTT*, thereby providing a first
84 solution to the learning dilemma for RSNNs. Furthermore *e-prop* also works for RSNNs
85 with more complex neuron models, such as LSNNs. This new learning paradigm for brain-
86 like network models elucidates how the brain could learn to recognize phonemes in spoken
87 language (Fig. 2), solve temporal credit assignment problems (Fig. 3), and acquire new
88 behaviors just from rewards (Fig. 4, 5).

89 In such reinforcement learning (RL) tasks the learner needs to explore its environment,
90 and find out which action gets rewarded in what state [13]. There is no “teacher” that
91 tells the learner what action would be optimal; in fact, the learner may never find that
92 out. Nevertheless learning methods such as *BPTT* are essential for a powerful form of
93 RL that is often referred to as Deep RL [14]. There one trains recurrent artificial neural
94 networks with internally generated teaching signals. We show here that Deep RL can in
95 principle also be carried out by neural networks of the brain, since *e-prop* approximates the
96 performance of *BPTT* also in this RL context. However another new ingredient is needed
97 to prove that. Previous work on Deep RL for solving complex tasks, such as winning
98 Atari games [14], required additional mechanisms to avoid well-known instabilities that
99 arise from using nonlinear function approximators, such as the use of several interacting
100 learners in parallel. Since this parallel learning scheme does not appear to be biologically
101 plausible, we introduce here a new method for avoiding learning instabilities: We show that
102 a suitable schedule for the lengths of learning episodes and learning rates also alleviates
103 learning instabilities in Deep RL.

104 We are not aware of previous work on online gradient descent learning methods for
105 RSNNs, neither for supervised learning nor for RL. There exists however preceding work
106 on online approximations of gradient descent for non-spiking neural networks based on
107 [15], which we review in the Discussion section.

108 The previous lack of powerful learning methods for RSNNs also affected the develop-
109 ment and use of neuromorphic computing hardware, which aims at a drastic reduction
110 in the energy consumption of AI implementations. A substantial fraction of this neuro-
111 morphic hardware, such as SpiNNaker [16] or Intel’s Loihi chip [17], implements RSNNs
112 and aims at on-chip training of these RSNNs. Although it does not matter here whether
113 the learning algorithm is biologically plausible, the excessive storage and offline processing
114 demands of *BPTT* make this option unappealing for neuromorphic hardware. Hence there
115 also exists a learning dilemma for RSNNs in neuromorphic hardware, which can be solved

116 with *e-prop*.

117 Results

118 Mathematical basis for *e-prop*

119 Spikes are modeled as binary variables z_j^t that assume value 1 if neuron j fires at time t ,
120 otherwise value 0. It is common in models to let t vary over small discrete time steps, e.g.
121 of 1 ms length. The goal of network learning is to find synaptic weights W that minimize
122 a given loss function E . E may depend on all or a subset of the spikes in the network.
123 E measures in the case of regression or classification learning the deviation of the actual
124 output y_k^t of each output neuron k at time t from its given target value $y_k^{*,t}$ (Fig. 1a).
125 In reinforcement learning (RL), the goal is to optimize the behavior of an agent in order
126 to maximize obtained rewards. In this case, E measures deficiencies of the current agent
127 policy to collect rewards.

128 The gradient $\frac{dE}{dW_{ji}}$ for the weight W_{ji} of the synapse from neuron i to neuron j tells
129 us how this weight should be changed in order to reduce E . The key innovation is that
130 a rigorous proof (see Methods) shows that this gradient can be represented as a sum over
131 the time steps t of the RSNN computation, where the second factor is just a local gradient
132 that does not depend on E :

$$\frac{dE}{dW_{ji}} = \sum_t \frac{dE}{dz_j^t} \cdot \left[\frac{dz_j^t}{dW_{ji}} \right]_{\text{local}} . \quad (1)$$

133 This local gradient can be represented as a sum of products of partial derivatives concern-
134 ing the hidden state of neuron j up to time t (equation (13)), which can be updated during
135 the forward computation of the RNN by a simple recursion (equation (14)). This term
136 $\left[\frac{dz_j^t}{dW_{ji}} \right]_{\text{local}}$ is not an approximation. Rather, it collects the maximal amount of informa-
137 tion about the network gradient $\frac{dE}{dW_{ji}}$ that can be computed locally in a forward manner.
138 Therefore it is the key-factor of *e-prop*. Since it reduces for simple neuron models – whose
139 internal state is fully captured by its membrane potential – to a variation of terms that
140 are commonly referred to as eligibility traces for synaptic plasticity [9], we also refer to

$$e_{ji}^t \stackrel{\text{def}}{=} \left[\frac{dz_j^t}{dW_{ji}} \right]_{\text{local}} \quad (2)$$

141 as eligibility trace. But most biological neurons have additional hidden variables that
142 change on a slower time scale, such as for example the firing threshold of a neuron with
143 firing threshold adaptation. Furthermore these slower processes in neurons are essential
144 for attaining with spiking neurons similarly powerful computing capabilities as LSTM
145 networks [3]. Hence the form that this eligibility trace e_{ji}^t takes for adapting neurons
146 (see equation (25)) is essential for understanding *e-prop*, and it is the main driver behind

147 the resulting qualitative jump in computing capabilities of RSNNs which are attainable
 148 through biologically plausible learning. Equations (1) and (2) yield the representation

$$\frac{dE}{dW_{ji}} = \sum_t L_j^t e_{ji}^t \quad (3)$$

149 of the loss gradient, where we refer to $L_j^t \stackrel{\text{def}}{=} \frac{dE}{dz_j^t}$ as the learning signal for neuron j . This
 150 equation defines a clear program for approximating the network loss gradient through local
 151 rules for synaptic plasticity: Change each weight W_{ji} at step t proportionally to $-L_j^t e_{ji}^t$, or
 152 accumulate these “tags” in a hidden variable that is translated occasionally into an actual
 153 weight change. Hence *e-prop* is an online learning method in a strict sense (see Fig. 1d
 154 and Movie S3). In particular, there is no need to unroll the network as for *BPTT*.

155 Since the ideal value $\frac{dE}{dz_j^t}$ of the learning signal L_j^t also captures influences which the
 156 current spike output z_j^t of neuron j may have on E via future spikes of other neurons, its
 157 precise value is in general not available at time t . We replace it by an approximation, such
 158 as $\frac{\partial E}{\partial z_j^t}$, which ignores these indirect influences. This approximation takes only currently
 159 arising losses at the output neurons k of the RSNN into account, and routes them with
 160 neuron-specific weights B_{jk} to the network neurons j (see Fig. 2a):

$$L_j^t = \sum_k B_{jk} \underbrace{(y_k^t - y_k^{*,t})}_{\text{deviation of output } k \text{ at time } t} . \quad (4)$$

161 Although this approximate learning signal L_j^t only captures errors that arise at the current
 162 time step t , it is combined in equation (3) with an eligibility trace e_{ji}^t that may reach
 163 far back into the past of neuron j (see Fig. 3b), thereby alleviating the need to solve
 164 the temporal credit assignment problem by propagating signals backwards in time (like in
 165 *BPTT*).

166 There are several strategies for choosing the weights B_{jk} for this online learning signal.
 167 In *symmetric e-prop* we set it equal to the corresponding weight W_{kj}^{out} of the synaptic
 168 connection from neuron j to output neuron k , as demanded by $\frac{\partial E}{\partial z_j^t}$. Note that this learning
 169 signal would actually implement $\frac{dE}{dz_j^t}$ exactly in the absence of recurrent connections in the
 170 network. Biologically more plausible are two variants of *e-prop* that avoid weight sharing:
 171 In *random e-prop* the values of all weights B_{jk} – even for neurons j that are not synaptically
 172 connected to output neuron k – are randomly chosen and remain fixed, similar to Broadcast
 173 Alignment for feedforward networks [18, 19, 20]. In *adaptive e-prop* we let in addition B_{jk}
 174 for neurons j that are synaptically connected to output neuron k evolve through a simple
 175 local plasticity rule that mirrors the plasticity rule applied to W_{kj}^{out} (see section S2.3).

176 Resulting synaptic plasticity rules (see Methods) look similar to previously proposed
 177 plasticity rules [9] for the special case of LIF neurons without slowly changing hidden
 178 variables. In particular they involve postsynaptic depolarization as one of the factors,
 179 similarly as the data-based Clopath-rule in [21], see section S6.4 in the supplement for an
 180 analysis.

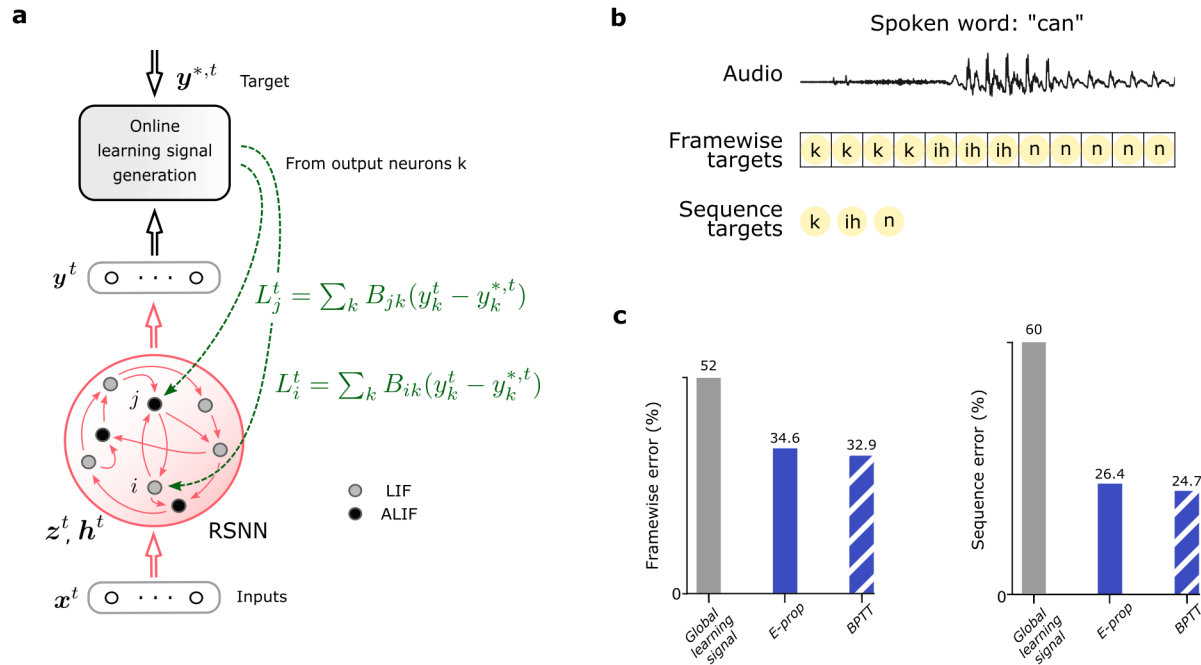


Figure 2: Comparison of the performance of *BPTT* and *e-prop* for learning phoneme recognition (TIMIT data set). **a**) Network architecture for *e-prop*, illustrated for an LSNN consisting of LIF and ALIF neurons. **b**) Input and target output for the two versions of TIMIT. **c**) Performance of *BPTT* and *symmetric e-prop* for LSNNs consisting of 800 neurons for framewise targets and 2400 for sequence targets (*random* and *adaptive e-prop* produced similar results, see Fig. S2). To obtain the “Global learning signal” baselines, the neuron specific feedbacks are replaced with global ones.

181 Comparing the performance of *e-prop* and *BPTT* for learning 182 spoken phoneme recognition

183 The phoneme recognition task TIMIT [22] is one of the most commonly used benchmarks
184 for temporal processing capabilities of different types of recurrent neural networks and
185 different learning approaches [23]. It comes in two versions. Both use, as input, acoustic
186 speech signals from sentences that are spoken by 630 speakers from 8 dialect regions of the
187 USA (see the top of Fig. 2b for a sample segment). In the simpler version, used for example
188 in [23], the goal is to recognize which of 61 phonemes is spoken in each 10 ms time frame
189 (“frame-wise classification”). In the more sophisticated version from [24], which achieved
190 an essential step toward human-level performance in speech-to-text transcription, the goal
191 is to recognize the sequence of phonemes in the entire spoken sentence independently
192 of their timing (“sequence transcription”). RSNNs consisting of LIF neurons do not even
193 reach with *BPTT* good performance for TIMIT [3]. Hence we are considering here LSNNs,
194 where a random subset of the neurons is a variation of the LIF model with firing rate
195 adaptation (ALIF neurons), see Methods. The name LSNN is motivated by the fact that

196 this special case of the RSNN model can achieve through training with *BPTT* similar
197 performance as an LSTM network [3].

198 *E-prop* approximates the performance of *BPTT* on LSNNs for both versions of TIMIT
199 very well, as shown in Fig. 2c. Furthermore LSNNs could solve the frame-wise classification
200 task without any neuron firing more frequently than 12 Hz (spike count taken over 32
201 spoken sentences), demonstrating that they operate in an energy efficient spike-coding –
202 rather than a rate-coding – regime. For the more difficult version of TIMIT we trained as
203 in [24] a complex LSNN consisting of a feedforward sequence of three recurrent networks.
204 Our results show that *e-prop* can also handle learning for such more complex network
205 structures very well. In Fig. S4 we show for comparison also the performance of *e-prop*
206 and *BPTT* for LSTM networks on the same tasks. These data show that for both versions
207 of TIMIT the performance of *e-prop* for LSNNs comes rather close to that of *BPTT* for
208 LSTM networks. In addition, they show that *e-prop* provides also for LSTM networks a
209 functionally powerful online learning method.

210 ***E-prop* performance for learning a task where temporal credit** 211 **assignment is difficult**

212 A hallmark of cognitive computations in the brain is the capability to go beyond a purely
213 reactive mode: to integrate diverse sensory cues over time, and to wait until the right
214 moment arrives for an action. A large number of experiments in neuroscience analyze
215 neural coding after learning such tasks (see e.g. [25, 11]). But it had remained unknown
216 how one can model the learning of such cognitive computations in RSNNs of the brain. In
217 order to test whether *e-prop* can solve this problem, we considered the same task that was
218 studied in the experiments of [25] and [11]. There a rodent moved along a linear track in
219 a virtual environment, where it encountered several visual cues on the left and right, see
220 Fig. 3a and Movie S1. Later, when it arrived at a T-junction, it had to decide whether to
221 turn left or right. It was rewarded when it turned to that side from which it had previously
222 received the majority of visual cues. This task is not easy to learn since the subject needs
223 to find out that it does not matter on which side the last cue was, or in which order the
224 cues were presented. Instead, the subject has to learn to count cues separately for each
225 side and to compare the two resulting numbers. Furthermore the cues need to be processed
226 properly long before a reward is given. We show in Fig. S5 that LSNNs can learn this task
227 via *e-prop* in exactly the same way just from rewards. But since the way how *e-prop* solves
228 the underlying temporal credit assignment problem is easier to explain for the supervised
229 learning version of this task, we discuss here the case where a teacher tells the subject at
230 the end of each trial what would have been the right decision. This still yields a challenging
231 scenario for any online learning method since non-zero learning signals L_j^t arise only during
232 the last 150 ms of a trial (Fig. 3b). Hence all synaptic plasticity has to take place during
233 these last 150 ms, long after the input cues have been processed. Nevertheless, *e-prop* is
234 able to solve this learning problem, see Fig. 3c and Movie S3. It just needs a bit more time
235 to reach the same performance level as offline learning via *BPTT* (see Movie S2). Whereas

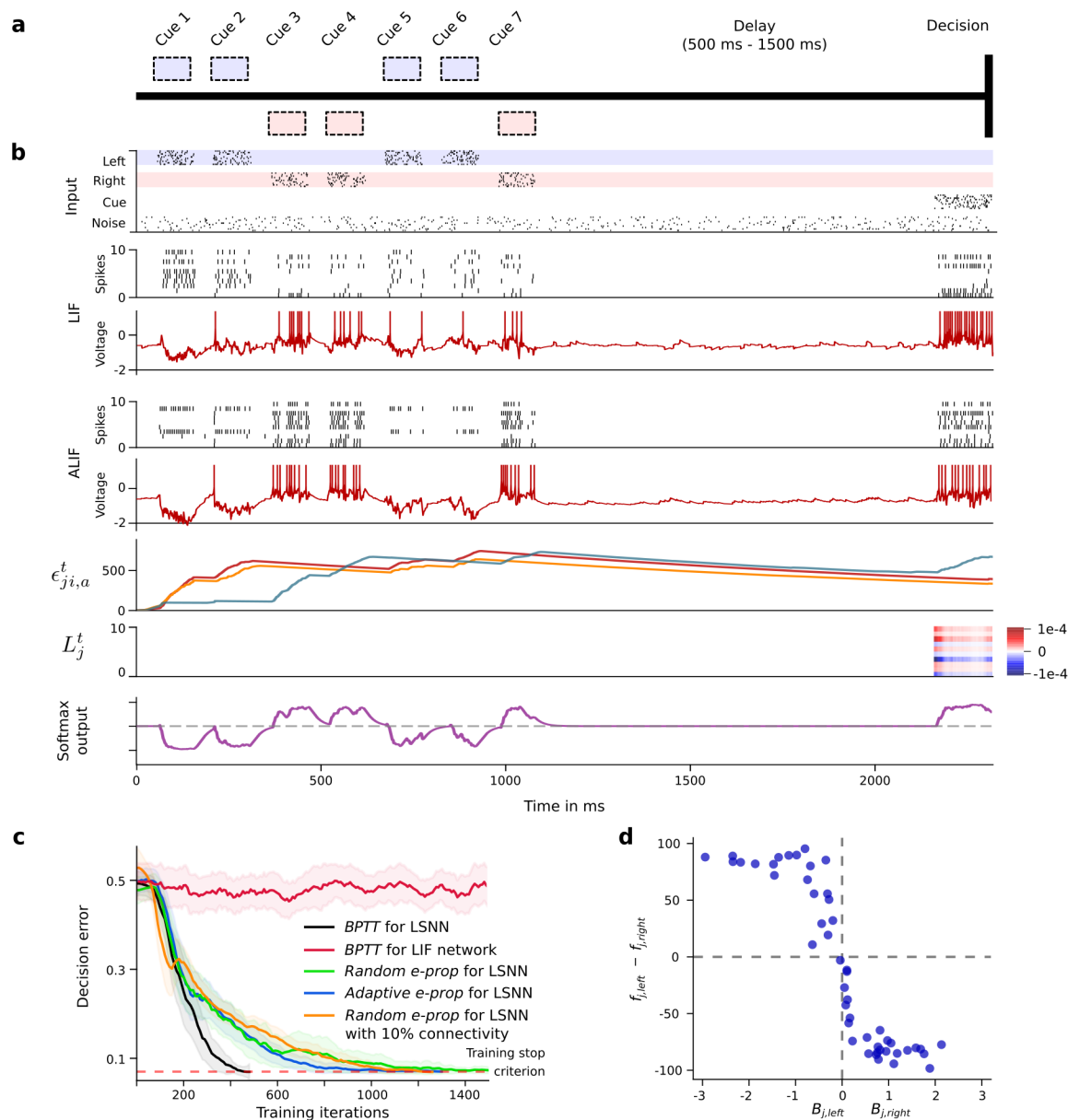


Figure 3: **Solving a task with difficult temporal credit assignment.** **a)** Setup of corresponding rodent experiments of [25] and [11], see Movie S1. **b)** Input spikes, spiking activity of 10 out of 50 sample LIF neurons and 10 out of 50 sample ALIF neurons, membrane potentials (more precisely: $v_j^t - A_j^t$) for two sample neurons j , 3 samples of slow components of eligibility traces, sample learning signals for 10 neurons and softmax network output. **c)** Learning curves for *BPTT* and two *e-prop* versions applied to LSNNs, and *BPTT* applied to an RSNN without adapting neurons (red curve). Orange curve shows learning performance of *e-prop* for a sparsely connected LSNN consisting of excitatory and inhibitory neurons (Dale’s law obeyed). The shaded areas are the 95%-confidence intervals of the mean accuracy computed with 20 runs. **d)** Correlation between the randomly drawn broadcast weights B_{jk} for $k = \text{left/right}$ for learning signals in *random e-prop* and resulting sensitivity to “left” and “right” input components after learning. $f_{j,\text{left}}$ ($f_{j,\text{right}}$) was the resulting average firing rate of neuron j during presentation of left (right) cues after learning.

236 this task can not even be solved by *BPTT* with a regular RSNN that has no adapting
237 neurons (red curve in Fig. 3c), all 3 previously discussed variations of *e-prop* can solve it
238 if the RSNN contains adapting neurons. We explain in section S2.5 how this task can also
239 be solved by sparsely connected LSNNs consisting of excitatory and inhibitory neurons:
240 by integrating stochastic rewiring [26] into *e-prop*.

241 But how can the neurons in the LSNN learn to record and count the input cues if all
242 the learning signals are identically 0 until the last 150 ms of a 2250 ms long trial (see 2nd
243 to last row of Fig. 3b)? For answering this question one should note that firing of a neuron
244 j at time t can affect the loss function E at a later time point $t' > t$ in two different ways:

- 245 i) By affecting future values of slow hidden variables of neuron j (e.g., its firing thresh-
246 old), which may then affect the firing of neuron j at t' , which in turn may directly
247 affect the loss function at time t' .
- 248 ii) By affecting the firing of other neurons j' at t' , which directly affects the loss function
249 at time t' .

250 In *symmetric* and *adaptive e-prop* one uses the partial derivative $\frac{\partial E}{\partial z_j^t}$ as learning signal
251 L_j^t for *e-prop* – instead of the online not available total derivative $\frac{dE}{dz_j^t}$. This blocks the
252 flow of gradient information along route ii. But the eligibility trace keeps the flow along
253 route i open. Therefore even *symmetric* and *adaptive e-prop* can solve the temporal credit
254 assignment problem of Fig. 3 through online learning: The gradient information that flows
255 along route i enables neurons to learn how to process the sensory cues at time points t
256 during the first 1050 ms, although this can affect the loss only at time points $t' > 2100$ ms
257 when the loss becomes non-zero. This is illustrated in the 3rd last row of Fig. 3b: The
258 slow component $\epsilon_{ji,a}^t$ of the eligibility traces e_{ji} of adapting neurons j decays with the
259 typical long time constant of firing rate adaptation (see equation (24) and Movie S3).
260 Since these traces stretch from the beginning of the trial into its last phase, they enable
261 learning of differential responses to “left” and “right” input cues that arrived over 1050 ms
262 before any learning signals become non-zero, as shown in the 2nd to last row of Fig. 3b.
263 Hence eligibility traces provide “highways into the future” for the propagation of gradient
264 information. These can be seen as biologically realistic replacements for the highways into
265 the past that *BPTT* employs during its backwards pass (see Movie S2).

266 This analysis also tells us when *symmetric e-prop* is likely to fail to approximate the
267 performance of *BPTT*: If forward propagation of gradient information cannot reach along
268 route i those later time points t' when the value of the loss function becomes salient. One
269 can artificially induce this in the experiment of Fig. 3 by adding to the LSNN – which has
270 the standard architecture shown in Fig. 2a – hidden layers of a feedforward SNN through
271 which the communication between the LSNN and the readout neurons has to flow. The
272 neurons j' of these hidden layers block route i, while leaving route ii open. Hence the task
273 of Fig. 3 can still be learnt with this modified network architecture by *BPTT*, but not by
274 *symmetric e-prop*, see Fig. S8.

275 Identifying tasks where the performance of *random e-prop* stays far behind that of
276 *BPTT* is more difficult, since error signals are sent there also to neurons that have no direct

277 connections to readout neurons. For deep feedforward networks it has been shown in [27]
278 that Broadcast Alignment, as defined in [20, 19], cannot reach the performance of Backprop
279 for difficult image classification tasks. Hence we expect that *random e-prop* will exhibit
280 corresponding deficiencies for difficult classification tasks with deep feedforward SNNs.
281 We are not aware of corresponding demonstrations of failures of Broadcast Alignment for
282 artificial RNNs, although they are likely to exist. Once they are found, they will probably
283 point to tasks where *random e-prop* fails for RSNNs. Currently we are not aware of any.

284 Fig. 3d provides insight into the functional role of the randomly drawn broadcast
285 weights in *random e-prop*: The difference of these weights determines for each neuron
286 j whether it learns to respond in the first phase of a trial more to cues from the left or
287 right. This observation suggests that neuron-specific learning signals for RSNNs have the
288 advantage that they can create a diversity of feature detectors for task-relevant network
289 inputs. Hence a suitable weighted sum of these feature detectors is later able to cancel
290 remaining errors at the network output, similarly as in the case of feedforward networks
291 [18].

292 We would like to point out that the use of the familiar actor-critic method in *reward-*
293 *based e-prop*, which we will discuss in the next section, provides an additional channel by
294 which information about future losses can gate synaptic plasticity of the *e-prop* learner at
295 the current time step t : Through the estimate $V(t)$ of the value of the current state, that
296 is simultaneously learnt via internally generated reward-prediction errors.

297 ***Reward-based e-prop***

298 Deep RL has significantly advanced the state of the art in machine learning and AI through
299 clever applications of *BPTT* to RL [14]. We found that one of the arguably most powerful
300 RL methods within the range of deep RL approaches that are not directly biologically
301 implausible, policy gradient in combination with actor-critic, can be implemented with
302 *e-prop*. This yields the biologically plausible and hardware friendly deep RL algorithm
303 *reward-based e-prop*. The LSNN learns here both an approximation to the value function
304 (the “critic”) and a stochastic policy (the “actor”). Neuron-specific learning signals are
305 combined in *reward-based e-prop* with a global signal that transmits reward prediction
306 errors (Fig. 4b). In contrast to the supervised case, where the learning signals L_j^t depend
307 on the deviation from an external target signal, the learning signals communicate here how
308 a stochastically chosen action deviates from the action mean that is currently proposed by
309 the network.

310 The resulting online synaptic plasticity rule (5) for deep RL is similar to equation (3),
311 except that a fading memory filter \mathcal{F}_γ is applied here to the term $L_j^t \bar{e}_{ji}^t$, where γ is the given
312 discount factor for future rewards and \bar{e}_{ji}^t denotes a low-pass filtered copy of the eligibility
313 trace e_{ji}^t (see Methods). This term is multiplied in the synaptic plasticity rule with the
314 reward prediction error $\delta^t = r^t + \gamma V^{t+1} - V^t$, where r^t is the reward received at time t .
315 This yields an instantaneous weight change of the form:

$$\Delta W_{ji}^t = -\eta \delta^t \mathcal{F}_\gamma(L_j^t \bar{e}_{ji}^t). \quad (5)$$

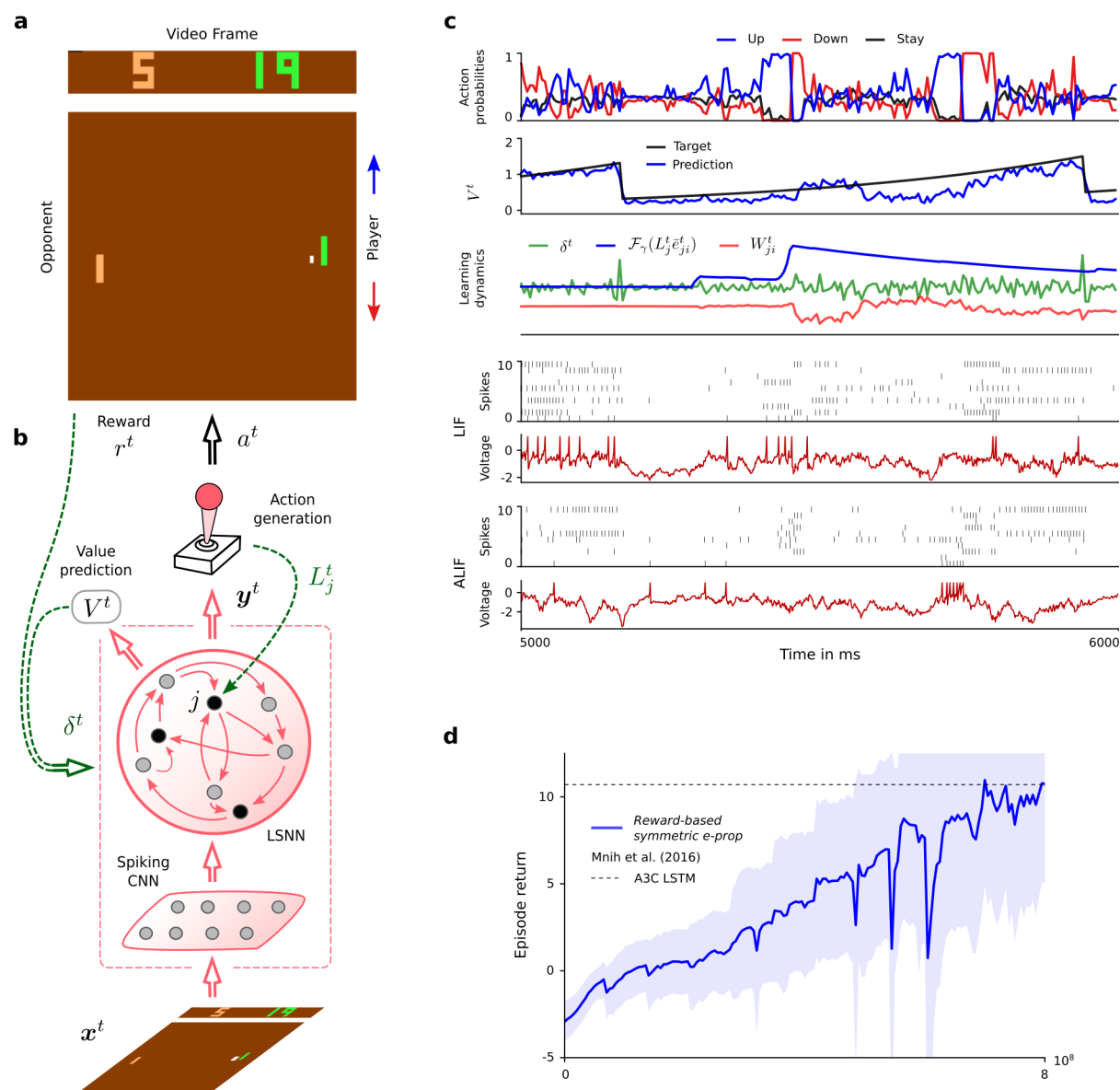


Figure 4: **Application of *e-prop* to the Atari game Pong.** **a)** Here the player (green paddle) has to outplay the opponent (light brown). A reward is acquired when the opponent cannot bounce back the ball (a small white square). To achieve this, the agent has to learn to hit the ball also with the edges of his paddle, which causes a less predictable trajectory. **b)** The agent is realized by an LSNN. The pixels of the current video frame of the game are provided as input. During processing of the stream of video frames by the LSNN, actions are generated by the stochastic policy in an online manner. At the same time, future rewards are predicted. The current error in prediction is fed back both to the LSNN and the spiking CNN that preprocesses the frames. (caption continued on next page)

Figure 4: (continued caption) **c**) Sample trial of the LSNN after learning with *reward-based e-prop*. From top to bottom: probabilities of stochastic actions, prediction of future rewards, learning dynamics of a random synapse (arbitrary units), spiking activity of 10 out of 240 sample LIF neurons and 10 out of 160 sample ALIF neurons, and membrane potentials (more precisely: $v_j^t - A_j^t$) for the two sample neurons j at the bottom of the spike raster above. **d**) Learning progress of the LSNN trained with *reward-based e-prop*, reported as the sum of collected rewards during an episode. The learning curve is averaged over 5 different runs and the shaded area represents the standard deviation. More information about the comparison between our results and A3C are given in section S5.3.

316 Previous 3-factor learning rules for RL were usually of the form $\Delta W^t = \eta \delta^t \bar{e}_{ji}^t$ [28, 9].
317 Hence they estimated gradients of the policy just by correlating the output of network
318 neurons with the reward prediction error. The learning power of this approach is known to
319 be quite limited due to high noise in the resulting gradient estimates. In contrast, in the
320 plasticity rule (5) for *reward-based e-prop* the eligibility traces are first combined with a
321 neuron specific feedback L_j^t , before they are multiplied with the reward prediction error δ^t .
322 We show in Methods analytically that this yields estimates of policy- and value gradients
323 similarly as in deep RL with *BPTT*. Furthermore, in contrast to previously proposed 3-
324 factor learning rules, this rule (5) is also applicable to LSNNs.

325 We tested *reward-based e-prop* on a classical benchmark task [14] for learning intelligent
326 behavior from rewards: Winning Atari video games provided by the Arcade Learning
327 Environment [29]. To win such game, the agent needs to learn to extract salient information
328 from the pixels of the game screen, and to infer the value of specific actions, even if rewards
329 are obtained in a distant future. In fact, learning to win Atari games is a serious challenge
330 for reinforcement learning even in machine learning [14]. Besides artificial neural networks
331 and *BPTT*, previous solutions also required experience replay (with a perfect memory of
332 many frames and action sequences that occurred much earlier) or an asynchronous training
333 of numerous parallel agents sharing synaptic weight updates. We show here that also an
334 LSNN can learn via *e-prop* to win Atari games, through online learning of a single agent.
335 This becomes possible with a single agent and without episode replay if the agent uses
336 a schedule of increasing episode lengths –with a learning rate that is inversely related to
337 that length. Using this scheme, an agent can experience diverse and uncorrelated short
338 episodes in the first phase of learning, producing useful skills. Subsequently, the agent can
339 fine-tune its policy using longer episodes.

340 First, we considered the well-known Atari game Pong (Fig. 4a). Here, the agent has to
341 learn to hit a ball in a clever way using up and down movements of his paddle. A reward
342 is obtained if the opponent cannot catch the ball. We trained an agent using *reward-based*
343 *e-prop* for this task, and show a sample trial in Fig. 4c and Movie S5. In contrast to
344 common deep RL solutions, the agent learns here in a strict online manner, receiving at
345 any time just the current frame of the game screen. In Panel d of Fig. 4 we demonstrate
346 that also this biologically realistic learning approach leads to a competitive score.

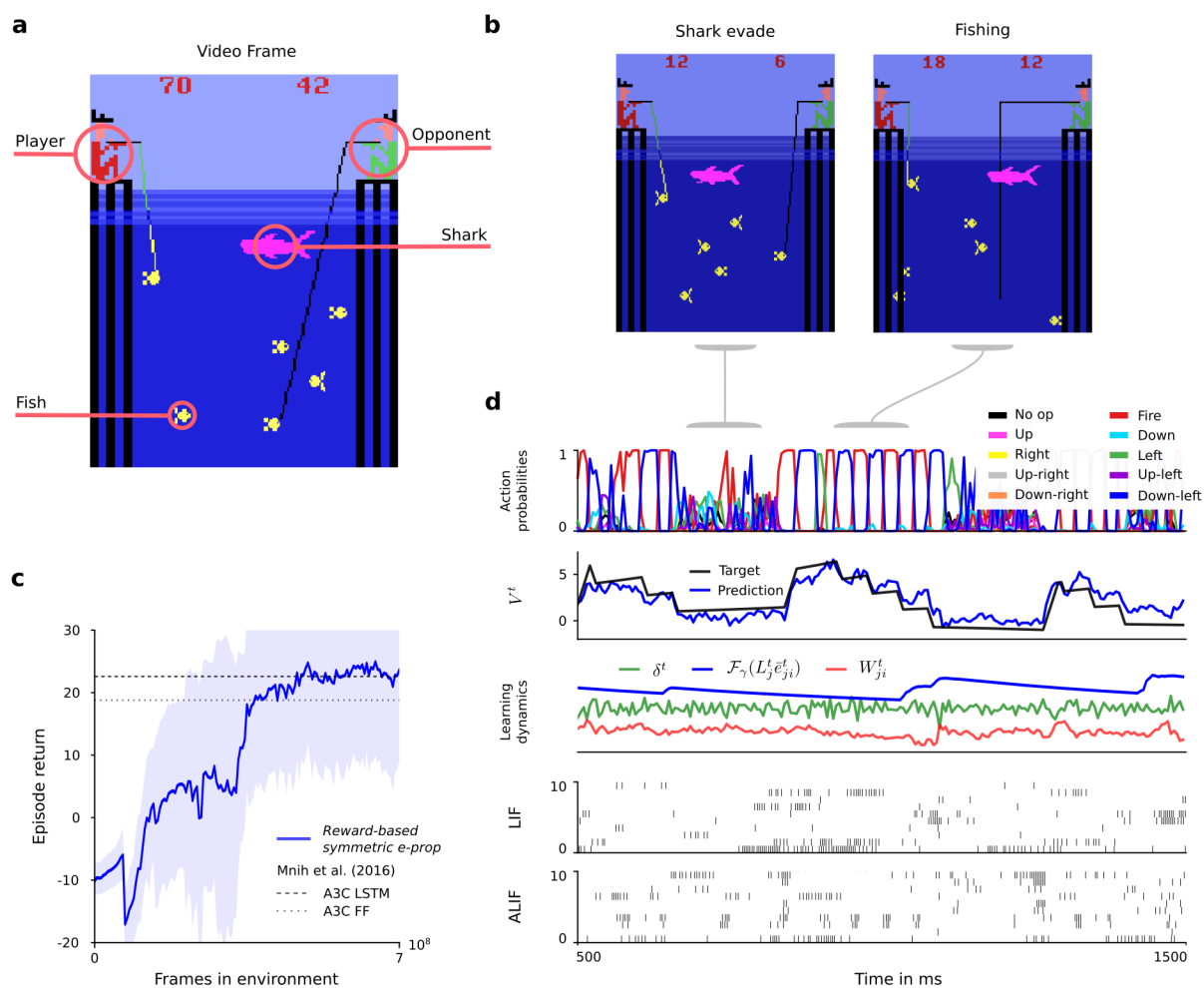


Figure 5: **Application of *e-prop* to learning to win the Atari game Fishing Derby.** **a)** Here the player has to compete against an opponent, and try to catch more fish from the sea. **b)** Once a fish has bit, the agent has to avoid that the fish gets touched by a shark. **c)** Sample trial of the trained network. From top to bottom: probabilities of stochastic actions, prediction of future rewards, learning dynamics of a random synapse (arbitrary units), spiking activity of 20 out of 180 sample LIF neurons and 20 out of 120 sample ALIF neurons. **d)** Learning curves of an LSNN trained with *reward-based e-prop* as in Fig. 4d.

347 If one does not insist on an online setting where the agent receives just the current
348 frame of the video screen but the last 4 frames, winning strategies for about half of the
349 Atari games can already be learnt by feedforward neural networks (see table S3 of [14]).
350 However, for other Atari games, such as Fishing Derby (Fig. 5a), it was even shown in [14]
351 that deep RL applied to LSTM networks achieves a substantially higher score than any deep
352 RL method for feedforward networks which was considered there. Hence, in order to test
353 the power of online *reward-based e-prop* also for those Atari games that require enhanced
354 temporal processing, we tested it on the Fishing Derby game. In this game, the agent has
355 to catch as many fish as possible while avoiding that the shark touches the fish with any
356 part of its body, and that the opponent catches the fish first. We show in Fig. 5c that
357 online *reward-based e-prop* applied to an LSNN does in fact reach the same performance
358 as reference offline algorithms applied to LSTM networks. We show a random trial after
359 learning in Fig. 5d, where we can identify two different learnt behaviors: 1) evading the
360 shark, 2) collecting fish. The agent has learnt to switch between these two behaviors as
361 required by the situation.

362 In general, we conjecture that variants of *reward-based e-prop* will be able to solve most
363 deep RL tasks that can be solved by online actor-critic methods in machine learning.

364 Discussion

365 We propose that in order to understand the computational function and neural coding of
366 neural networks in the brain, one needs to understand the organization of the plasticity
367 mechanisms that install and maintain these. So far *BPTT* was the only candidate for
368 that, since no other learning method provided sufficiently powerful computational function
369 to RSNN models. But since *BPTT* is not viewed to be biologically realistic [5], it does
370 not help us to understand learning in the brain. *E-prop* offers a solution to this dilemma,
371 since it does not require biologically unrealistic mechanisms, but still enables RSNNs to
372 learn difficult computational tasks, in fact almost as well as *BPTT*. Furthermore it enables
373 RSNNs to solve these tasks in an energy efficient sparse firing regime, rather than resorting
374 to rate coding.

375 *E-prop* relies on two types of signals that are abundantly available in the brain, but
376 whose precise role for learning have not yet been understood: eligibility traces and learning
377 signals. Since *e-prop* is based on a transparent mathematical principle (see equation (3)),
378 it provides a normative model for both types of signals, as well as for synaptic plasticity
379 rules. Interestingly, the resulting learning model suggests that a characteristic aspect of
380 many biological neurons – the presence of slowly changing hidden variables – provides a
381 possible solution to the problem how a RSNN can learn without error signals that propagate
382 backwards in time: Slowly changing hidden variables of neurons cause eligibility traces
383 that propagate forward over longer time spans, and are therefore able to coincide with
384 later arising instantaneous error signals (see Fig. 3b).

385 The theory of *e-prop* makes a concrete experimentally testable prediction: that the
386 time constant of the eligibility trace for a synapse is correlated with the time constant for

387 the history-dependence of the firing activity of the postsynaptic neuron. It also suggests
388 that the experimentally found diverse time constants of the firing activity of populations
389 of neurons in different brain areas [30] are correlated with their capability to handle corre-
390 sponding ranges of delays in temporal credit assignment for learning.

391 Finally, *e-prop* theory provides a hypothesis for the functional role of the experimentally
392 found diversity of dopamine signals to different populations of neurons [11]. Whereas
393 previous theories of reward-based learning required that the same learning signal is sent
394 to all neurons, the basic equation (1) for *e-prop* postulates that ideal top-down learning
395 signals to a population of neurons depend on its impact on the network performance (loss
396 function), and should therefore be target specific (see Fig. 2c and section S6.2). In fact, the
397 learning-to-learn result for *e-prop* in [31] suggests that prior knowledge about the possible
398 range of learning tasks for a brain area could optimize top-down learning signals even
399 further on an evolutionary time scale, thereby enabling for example learning from few or
400 even a single trial.

401 Previous methods for training RSNNs did not aim at approximating *BPTT*. Instead
402 some of them were relying on control theory to train a chaotic reservoir of spiking neurons
403 [32, 33, 34]. Others used the FORCE algorithm [35, 36] or variants of it [37, 38, 39,
404 35]. However the FORCE algorithm was not argued to be biologically realistic, since the
405 plasticity rule for each synaptic weight requires knowledge of the current values of all other
406 synaptic weights. The generic task considered in [35] was to learn with supervision how
407 to generate patterns. We show in Figs. S1, S7, and Movie S4 that RSNNs can learn such
408 tasks also with a biologically plausible learning method: *e-prop*.

409 Several methods for approximating stochastic gradient descent in *feedforward* networks
410 of spiking neurons have been proposed, see e.g. [40, 41, 42, 43, 44]. These employ –
411 like *e-prop* – a pseudo-gradient to overcome the non-differentiability of a spiking neuron,
412 as proposed previously in [45, 46]. [40, 42, 43] arrive at a synaptic plasticity rule for
413 feedforward networks that consists – like *e-prop* – of the product of a learning signal and
414 a derivative (eligibility trace) that describes the dependence of a spike of a neuron j on
415 the weight of an afferent synapse W_{ji} . But in a recurrent network the spike output of j
416 depends on W_{ji} also indirectly, via loops in the network that allow that a spike of neuron
417 j contributes to the firing of other neurons, which in turn affect firing of the presynaptic
418 neuron i . Hence the corresponding eligibility trace can no longer be locally computed if
419 one transfers these methods for feedforward networks to recurrently connected networks.
420 Therefore [40] suggests the need to investigate extensions of their approach to RSNNs.

421 Previous work on the design of online gradient descent learning algorithms for *non-*
422 *spiking* RNNs was based on real-time recurrent learning (RTRL) [15]. RTRL itself has
423 rarely been used since its computational complexity per time-step is $\mathcal{O}(n^4)$, if n is the
424 number of neurons. But interesting approximations to RTRL have subsequently been pro-
425 posed (see [47] for a review): some stochastic approximations [48] which are $\mathcal{O}(n^3)$ or only
426 applicable for small networks [49], and also recently two deterministic $\mathcal{O}(n^2)$ approxima-
427 tions [50, 51]. The latter were in fact written at the same time as the first publication
428 of *e-prop* [31]. A structural difference between this paper and [50] is that their approach
429 requires that learning signals are transmitted between the neurons in the RNN, with sepa-

430 rately learnt weights. [51] derived for rate based neurons a learning rule similar to *random*
431 *e-prop*. But this work did not address other forms of learning than supervised regression,
432 such as RL, nor learning in networks of spiking neurons, or in more powerful types of RNNs
433 with slow hidden variables such as LSTM networks or LSNNs.

434 *E-prop* also has complexity $\mathcal{O}(n^2)$, in fact $\mathcal{O}(S)$ if S is the number of synaptic connec-
435 tions. This bound is optimal -except for the constant factor- since this is the asymptotic
436 complexity of just simulating the RNN. The key point of *e-prop* is that the general form (13)
437 of its eligibility trace collects all contributions to the loss gradient that can be locally com-
438 puted in a feedforward manner. This general form enables applications to spiking neurons
439 with slowly varying hidden variables, such as neurons with firing rate adaptation, which
440 are essential ingredients of RSNNs to reach the computational power of LSTM networks
441 [3]. We believe that this approach can be extended in future work –with a suitable choice
442 of pseudo-derivatives—to a wide range of biologically more realistic neuron models. It also
443 enables the combination of these rigorously derived eligibility traces with – semantically
444 identical but algorithmically very different – eligibility traces from RL for *reward-based*
445 *e-prop* (equation (5)), thereby bringing the power of deep RL to RSNNs. As a result, we
446 were able to show in Fig. 2 - 5 that RSNNs can learn with the biologically plausible rules
447 for synaptic plasticity that arise from the *e-prop* theory to solve tasks such as phoneme
448 recognition, integrating evidence over time and waiting for the right moment to act, and
449 winning Atari games. These are tasks that are fundamental for modern learning-based AI,
450 but have so far not been solved with RSNNs. Hence *e-prop* provides a new perspective of
451 the major open question how intelligent behavior can be learnt and controlled by neural
452 networks of the brain.

453 Apart from obvious consequences of *e-prop* for research in neuroscience and cognitive
454 science, *e-prop* also provides an interesting new tool for approaches in machine learning
455 where *BPTT* is replaced by approximations in order to improve computational efficiency.
456 We have already shown in Fig. S4 that *e-prop* provides a powerful online learning method
457 for LSTM networks. Furthermore, the combination of eligibility traces from *e-prop* with
458 synthetic gradients from [52] even improves performance of LSTM networks for difficult
459 machine learning problems such as the copy-repeat task and the Penn Treebank word
460 prediction task [31]. Other future extensions of *e-prop* could explore a combination with
461 attention-based models in order to cover multiple timescales.

462 Finally, *e-prop* suggests a promising new approach for realizing powerful on-chip learn-
463 ing of RSNNs on neuromorphic chips. Whereas *BPTT* is not within the reach of current
464 neuromorphic hardware, an implementation of *e-prop* appears to offer no serious hurdle.
465 Our results show that an implementation of *e-prop* will provide a qualitative jump in
466 on-chip learning capabilities of neuromorphic hardware.

467 Methods

468 Table of Contents:

- 469 • Network models
- 470 • Conventions
- 471 • Mathematical basis for *e-prop*
- 472 • Eligibility traces
- 473 • Derivation of eligibility traces for concrete neuron models
- 474 • Synaptic plasticity rules resulting from *e-prop*
- 475 • *Reward-based e-prop*: application of *e-prop* to deep RL

476 Network models

477 To exhibit the generality of the *e-prop* approach, we define the dynamics of recurrent neural
 478 networks using a general formalism that is applicable to many recurrent neural network
 479 models, not only to RSNNs and LSNNs. Also non-spiking models such as LSTM networks
 480 fit under this formalism (see section S4.3 in the supplement). The network dynamics is
 481 summarized by the computational graph in Fig. 6. It uses the function M to define the
 482 update of the hidden state of a neuron j : $\mathbf{h}_j^t = M(\mathbf{h}_j^{t-1}, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$, where \mathbf{W}_j gathers
 483 the weights of synapses arriving at neuron j , and f to define the update of the observable
 484 state of a neuron j : $z_j^t = f(\mathbf{h}_j^t, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$ (f simplifies to $z_j^t = f(\mathbf{h}_j^t)$ for LIF and ALIF
 485 neurons). We chose a discrete time step of $\delta t = 1$ ms for all our simulations. Control
 486 experiments with smaller time steps for the task of Fig. 3, reported in Fig. S6, suggest that
 487 the size of the time step has no significant impact on the performance of *e-prop*.

488 **LIF neurons.** Each LIF neuron has a one dimensional internal state – or hidden variable
 489 – h_j^t that consists only of the membrane potential v_j^t . The observable state $z_j^t \in \{0, 1\}$ is
 490 binary, indicating a spike ($z_j^t = 1$) or no spike ($z_j^t = 0$) at time t . The dynamics of the LIF
 491 model is defined by the equations:

$$v_j^{t+1} = \alpha v_j^t + \sum_{i \neq j} W_{ji}^{\text{rec}} z_i^t + \sum_i W_{ji}^{\text{in}} x_i^{t+1} - z_j^t v_{\text{th}} \quad (6)$$

$$z_j^t = H(v_j^t - v_{\text{th}}). \quad (7)$$

492 W_{ji}^{rec} (W_{ji}^{in}) is the synaptic weight from network (input) neuron i to neuron j . The decay
 493 factor α in (6) is given by $e^{-\delta t/\tau_m}$, where τ_m (typically 20 ms) is the membrane time
 494 constant. δt denotes the discrete time step size, which is set to 1 ms in our simulations. H
 495 denotes the Heaviside step function. Note that we deleted in equation (6) the factor $1 - \alpha$
 496 that occurred in the corresponding equation (4) in the supplement of [3]. This simplifies
 497 the notation in our derivations, and has no impact on the model if parameters like the
 498 threshold voltage are scaled accordingly.

499 Due to the term $-z_j^t v_{\text{th}}$ in equation (6), the neurons membrane potential is reduced
 500 by a constant value after an output spike, which relates our model to the spike response
 501 model [53]. To introduce a simple model of neuronal refractoriness, we further assume
 502 that z_j^t is fixed to 0 after each spike of neuron j for a short refractory period of 2 to 5 ms
 503 depending on the simulation.

504 **LSNNs.** According to the database of the Allen Institute [2] a fraction of neurons be-
 505 tween roughly 20 % (in mouse visual cortex) and 40 % (in the human frontal lobe) exhibit
 506 spike frequency adaptation (SFA). It had been shown in [3] that the inclusion of neuron
 507 models with SFA – via a time-varying firing threshold as slow hidden variable – drastically
 508 enhances computing capabilities of RSNN models. Hence we consider here the same simple
 509 model for neurons with SFA as in [3], to which we refer as adaptive LIF (ALIF) neuron.
 510 This model is basically the same as the GLIF₂ model in the Technical White paper on
 511 generalized LIF (GLIF) models from [2]. LSNNs are recurrently connected networks that
 512 consist of LIF and ALIF neurons. ALIF neurons j have a second hidden variable a_j^t , which
 513 denotes the variable component of its firing threshold. As a result, their internal state is
 514 a 2 dimensional vector $\mathbf{h}_j^t \stackrel{\text{def}}{=} [v_j^t, a_j^t]$. Their threshold potential A_j^t increases with every
 515 output spike and decreases exponentially back to the baseline threshold v_{th} . This can be
 516 described by

$$A_j^t = v_{\text{th}} + \beta a_j^t, \quad (8)$$

$$z_j^t = H(v_j^t - A_j^t), \quad (9)$$

517 with a threshold adaptation according to

$$a_j^{t+1} = \rho a_j^t + z_j^t, \quad (10)$$

518 where the decay factor ρ is given by $e^{-\delta t/\tau_a}$, and τ_a is the adaptation time constant that is
 519 typically chosen to be in the range of the time span of the length of the working memory
 520 that is a relevant for a given task. This is a very simple model for a neuron with spike
 521 frequency adaptation [3]. We refer to [53, 54, 55] for experimental data and other neuron
 522 models. We refer to a recurrent network of spiking neurons (RSNN) as LSNN if some of
 523 its neurons are adaptive. We chose a fraction between 25 and 40 % of the neurons to be
 524 adapting, like in the data from neocortex [2], with time constants that are roughly on the
 525 same time scale as the tasks for which the network is trained.

526 In relation to the more general formalism represented in the computational graph in
 527 Fig. 6, equations (6) and (10) define $M(\mathbf{h}_j^{t-1}, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$, and equations (7) and (9)
 528 define $f(\mathbf{h}_j^t)$.

529 **Gradient descent for RSNNs.** Gradient descent is problematic for spiking neurons
 530 because of the step function H in equation (7). We overcome this issue as in [56, 3]: The
 531 non-existing derivative $\frac{\partial z_j^t}{\partial v_j^t}$ is replaced in simulations by a simple nonlinear function of the
 532 membrane potential that is called the pseudo-derivative. Outside of the refractory period,

533 we choose a pseudo-derivative of the form $\psi_j^t = \frac{1}{v_{th}} \gamma_{pd} \max\left(0, 1 - \left| \frac{v_j^t - A_j^t}{v_{th}} \right| \right)$ where $\gamma_{pd} = 0.3$
 534 for ALIF neurons, and for LIF neurons A_j^t is replaced by v_{th} . During the refractory period
 535 the pseudo derivative is set to 0.

536 **Network output and loss functions.** We assume that network outputs y_k^t are real-
 537 valued and produced by leaky output neurons (readouts) k , which are not recurrently
 538 connected:

$$y_k^t = \kappa y_k^{t-1} + \sum_j W_{kj}^{out} z_j^t + b_k^{out}, \quad (11)$$

539 where $\kappa \in [0, 1]$ defines the leak and b_k^{out} denotes the output bias. The leak factor κ is
 540 given for spiking neurons by $e^{-\delta t / \tau_{out}}$, where τ_{out} is the membrane time constant. Note that
 541 for non-spiking neural networks (such as for LSTM networks), temporal smoothing of the
 542 network observable state is not necessary. In this case, one can use $\kappa = 0$.

543 The loss function $E(\mathbf{z}^1, \dots, \mathbf{z}^T)$ quantifies the network performance. We assume that
 544 it depends only on the observable states $\mathbf{z}^1, \dots, \mathbf{z}^T$ of the network neurons. For instance,
 545 for a regression problem we define E as the mean square error $E = \frac{1}{2} \sum_{t,k} (y_k^t - y_k^{*,t})^2$
 546 between the network outputs y_k^t and target values $y_k^{*,t}$. For classification or RL tasks the
 547 loss function E has to be re-defined accordingly.

548 Conventions

549 **Notation for derivatives.** We distinguish the total derivative $\frac{dE}{d\mathbf{z}^t}(\mathbf{z}^1, \dots, \mathbf{z}^T)$, which
 550 takes into account how E depends on \mathbf{z}^t also indirectly through influence of \mathbf{z}^t on the other
 551 variables $\mathbf{z}^{t+1}, \dots, \mathbf{z}^T$, and the partial derivative $\frac{\partial E}{\partial \mathbf{z}^t}(\mathbf{z}^1, \dots, \mathbf{z}^T)$ which quantifies only the
 552 direct dependence of E on \mathbf{z}^t .

553 Analogously for the hidden state $\mathbf{h}_j^t = M(\mathbf{h}_j^{t-1}, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$ the partial derivative
 554 $\frac{\partial M}{\partial \mathbf{h}_j^{t-1}}$ denotes the partial derivative of M with respect to \mathbf{h}_j^{t-1} . It only quantifies the
 555 direct influence of \mathbf{h}_j^{t-1} on \mathbf{h}_j^t and it does not take into account how \mathbf{h}_j^{t-1} indirectly in-
 556 fluences \mathbf{h}_j^t via the observable states \mathbf{z}^{t-1} . To improve readability we also use the follow-
 557 ing abbreviations: $\frac{\partial \mathbf{h}_j^t}{\partial \mathbf{h}_j^{t-1}} \stackrel{\text{def}}{=} \frac{\partial M}{\partial \mathbf{h}_j^{t-1}}(\mathbf{h}_j^{t-1}, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$, $\frac{\partial \mathbf{h}_j^t}{\partial W_{ji}} \stackrel{\text{def}}{=} \frac{\partial M}{\partial W_{ji}}(\mathbf{h}_j^{t-1}, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$, and
 558 $\frac{\partial z_j^t}{\partial \mathbf{h}_j^t} \stackrel{\text{def}}{=} \frac{\partial f}{\partial h^t}(\mathbf{h}_j^t, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$.

559 **Notation for temporal filters.** For ease of notation we use the operator \mathcal{F}_α to denote
 560 the low-pass filter such that, for any time series x_t :

$$\mathcal{F}_\alpha(x^t) = \alpha \mathcal{F}_\alpha(x^{t-1}) + x^t, \quad (12)$$

561 and $\mathcal{F}_\alpha(x^0) = x^0$. In the specific case of the time series z_j^t and e_{ji}^t , we simplify notation
 562 further and write \bar{z}_j^t and \bar{e}_{ji}^t for $\mathcal{F}_\alpha(z_j^t)$ and $\mathcal{F}_\alpha(e_{ji}^t)$.

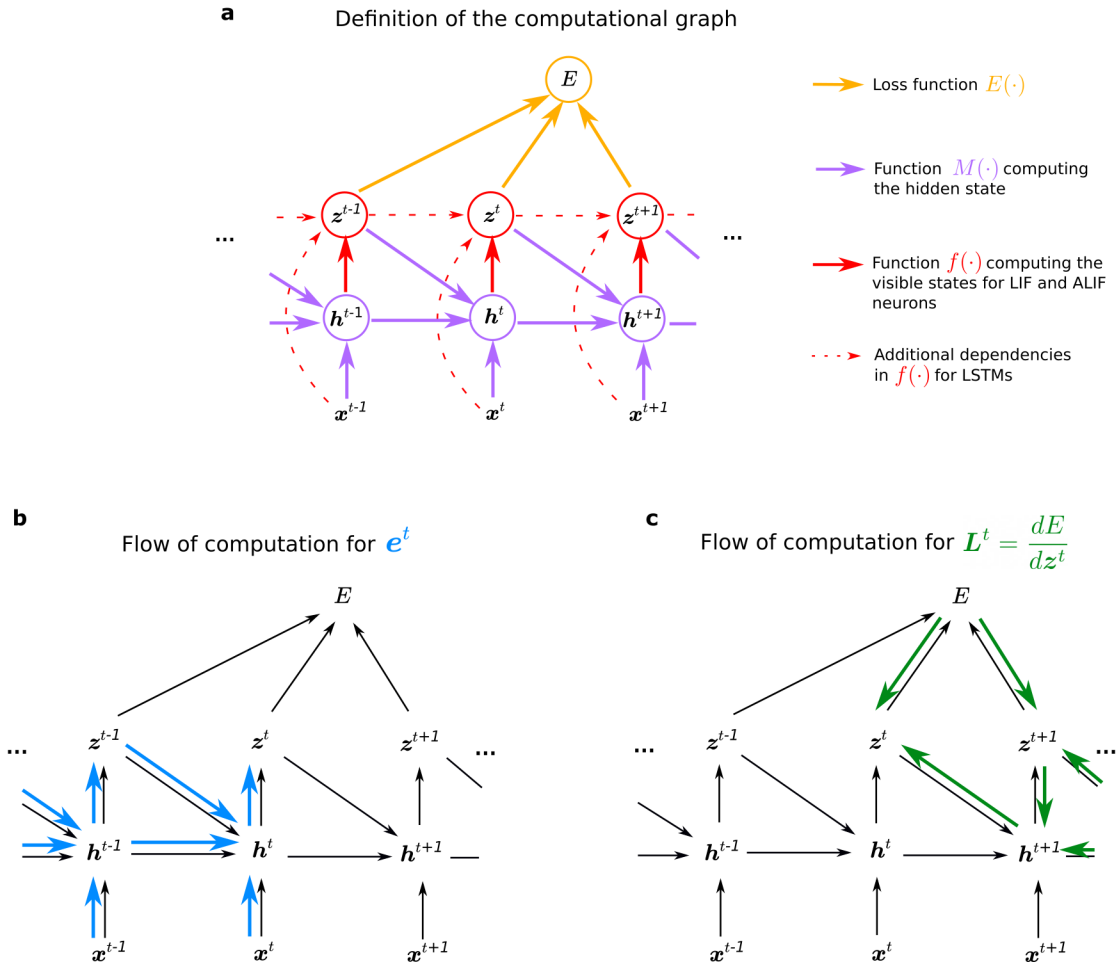


Figure 6: Computational graph and gradient propagations **a)** Assumed mathematical dependencies between hidden neuron states \mathbf{h}_j^t , neuron outputs \mathbf{z}^t , network inputs \mathbf{x}^t , and the loss function E through the mathematical functions $E(\cdot)$, $M(\cdot)$, $f(\cdot)$ are represented by coloured arrows. **b-c)** The flow of computation for the two components \mathbf{e}^t and \mathbf{L}^t that merge into the loss gradients of equation (3) can be represented in similar graphs. **b)** Following equation (14), the flow of the computation of the eligibility traces e_{ji}^t is going forward in time. **c)** Instead the ideal learning signals $L_j^t = \frac{dE}{dz_j^t}$ require to propagate gradients backward in time. Hence while e_{ji}^t is computed exactly, L_j^t is approximated in *e-prop* applications to yield an online learning algorithm.

563 Mathematical basis for *e-prop*

564 We provide here the proof of the fundamental equation (1) for *e-prop*

$$\frac{dE}{dW_{ji}} = \sum_t \frac{dE}{dz_j^t} \cdot \left[\frac{dz_j^t}{dW_{ji}} \right]_{\text{local}},$$

with the new eligibility trace

$$e_{ji}^t \stackrel{\text{def}}{=} \left[\frac{dz_j^t}{dW_{ji}} \right]_{\text{local}} \stackrel{\text{def}}{=} \frac{\partial z_j^t}{\partial \mathbf{h}_j^t} \underbrace{\sum_{t' \leq t} \frac{\partial \mathbf{h}_j^t}{\partial \mathbf{h}_j^{t-1}} \cdots \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}} \cdot \frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}}}_{\stackrel{\text{def}}{=} \boldsymbol{\epsilon}_{ji}^t}. \quad (13)$$

565 For spiking neurons j we replace the first factor $\frac{\partial z_j^t}{\partial \mathbf{h}_j^t}$ of e_{ji}^t by the pseudo-derivative, see
 566 [3, 4, 56]. The second factor $\boldsymbol{\epsilon}_{ji}^t$, which we call eligibility vector, obviously satisfies the
 567 recursive equation

$$\boldsymbol{\epsilon}_{ji}^t = \frac{\partial \mathbf{h}_j^t}{\partial \mathbf{h}_j^{t-1}} \cdot \boldsymbol{\epsilon}_{ji}^{t-1} + \frac{\partial \mathbf{h}_j^t}{\partial W_{ji}}, \quad (14)$$

568 where \cdot denotes the dot product. This provides the rule for the online computation of $\boldsymbol{\epsilon}_{ji}^t$,
 569 and hence of $e_{ji}^t = \frac{\partial z_j^t}{\partial \mathbf{h}_j^t} \cdot \boldsymbol{\epsilon}_{ji}^t$.

570 We start from a classical factorization of the loss gradients in recurrent neural networks
 571 that arises for instance in equation (12) of [57] to describe *BPTT*. This classical factor-
 572 ization can be justified by unrolling an RNN into a large feedforward network where each
 573 layer (l) represents one time step. In a feedforward network the loss gradients with respect
 574 to the weights $W_{ji}^{(l)}$ of layer l are given by $\frac{dE}{dW_{ji}^{(l)}} = \frac{dE}{dh_j^{(l)}} \frac{\partial h_j^{(l)}}{\partial W_{ji}^{(l)}}$. But as the weights are shared
 575 across the layers when representing a recurrent network, the summation of these gradients
 576 over the layers l of the unrolled RNN yields this classical factorization of the loss gradients:

$$\frac{dE}{dW_{ji}} = \sum_{t'} \frac{dE}{d\mathbf{h}_j^{t'}} \cdot \frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}}. \quad (15)$$

577 Note that the first factor $\frac{dE}{d\mathbf{h}_j^{t'}}$ in these products also needs to take into account how the
 578 internal state \mathbf{h}_j of neuron j evolves during subsequent time steps, and whether it influences
 579 firing of j at later time steps. This is especially relevant for ALIF neurons and other
 580 biologically realistic neuron models with slowly changing internal states. Note that this
 581 first factor of (15) is replaced in the *e-prop* equation (13) by the derivative $\frac{dE}{dz_j^{t'}}$ of E with
 582 regard to the observable variable $z_j^{t'}$. There the evolution of the internal state of neuron j
 583 is pushed into the second factor, the eligibility trace e_{ji} , which collects in *e-prop* all online
 584 computable factors of the loss gradient that just involve neurons j and i .

585 Now we show that one can re-factorize the expression (15) and prove that the loss
 586 gradients can also be computed using the new factorization (13) that underlies *e-prop*. In
 587 the steps of the subsequent proof until equation (19), we decompose the term $\frac{dE}{d\mathbf{h}_j^t}$ into a
 588 series of learning signals $L_j^t = \frac{dE}{dz_j^t}$ and local factors $\frac{\partial \mathbf{h}_j^{t+1}}{\partial \mathbf{h}_j^t}$ for $t \geq t'$. Those local factors
 589 will later be used to transform the partial derivative $\frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}}$ from equation (15) into the
 590 eligibility vector $\boldsymbol{\epsilon}_{ji}^t$ that integrates the whole history of the synapse up to time t , not just
 591 a single time step. To do so, we express $\frac{dE}{d\mathbf{h}_j^{t'}}$ recursively as a function of the same derivative
 592 at the next time step $\frac{dE}{d\mathbf{h}_j^{t'+1}}$ by applying the chain rule at the node \mathbf{h}_j^t for $t = t'$ of the
 593 computational graph shown in Fig. 6c:

$$\frac{dE}{d\mathbf{h}_j^{t'}} = \frac{dE}{dz_j^{t'}} \frac{\partial z_j^{t'}}{\partial \mathbf{h}_j^{t'}} + \frac{dE}{d\mathbf{h}_j^{t'+1}} \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}} \quad (16)$$

$$= L_j^{t'} \frac{\partial z_j^{t'}}{\partial \mathbf{h}_j^{t'}} + \frac{dE}{d\mathbf{h}_j^{t'+1}} \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}}, \quad (17)$$

594 where we defined the learning signal $L_j^{t'}$ as $\frac{dE}{dz_j^{t'}}$. The resulting recursive expansion ends at
 595 the last time step T of the computation of the RNN, i.e., $\frac{dE}{d\mathbf{h}_j^{T+1}} = 0$. If one repeatedly sub-
 596 stitutes the recursive formula (17) into the classical factorization (15) of the loss gradients,
 597 one gets:

$$\frac{dE}{dW_{ji}} = \sum_{t'} \left(L_j^{t'} \frac{\partial z_j^{t'}}{\partial \mathbf{h}_j^{t'}} + \frac{dE}{d\mathbf{h}_j^{t'+1}} \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}} \right) \cdot \frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}} \quad (18)$$

$$= \sum_{t'} \left(L_j^{t'} \frac{\partial z_j^{t'}}{\partial \mathbf{h}_j^{t'}} + (L_j^{t'+1} \frac{\partial z_j^{t'+1}}{\partial \mathbf{h}_j^{t'+1}} + (\dots) \frac{\partial \mathbf{h}_j^{t'+2}}{\partial \mathbf{h}_j^{t'+1}}) \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}} \right) \cdot \frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}}. \quad (19)$$

598 The following equation is the main equation for understanding the transformation from
 599 *BPTT* into *e-prop*. The key idea is to collect all terms $\frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}}$ which are multiplied with
 600 the learning signal L_j^t at a given time t . These are only terms that concern events in
 601 the computation of neuron j up to time t , and they do not depend on other future losses
 602 or variable values. To this end, we write the term in parentheses in equation (19) into
 603 a second sum indexed by t and exchange the summation indices to pull out the learning
 604 signal L_j^t . This expresses the loss gradient of E as a sum of learning signals L_j^t multiplied
 605 by some factor indexed by ji , which we define as the eligibility trace $e_{ji}^t \in \mathbb{R}$. The main
 606 factor of it is the eligibility vector $\boldsymbol{\epsilon}_{ji}^t \in \mathbb{R}^d$, which has the same dimension as the hidden

607 state \mathbf{h}_j^t :

$$\frac{dE}{dW_{ji}} = \sum_{t'} \sum_{t \geq t'} L_j^t \frac{\partial z_j^t}{\partial \mathbf{h}_j^t} \frac{\partial \mathbf{h}_j^t}{\partial \mathbf{h}_j^{t-1}} \cdots \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}} \cdot \frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}} \quad (20)$$

$$= \sum_t L_j^t \frac{\partial z_j^t}{\partial \mathbf{h}_j^t} \underbrace{\sum_{t' \leq t} \frac{\partial \mathbf{h}_j^t}{\partial \mathbf{h}_j^{t-1}} \cdots \frac{\partial \mathbf{h}_j^{t'+1}}{\partial \mathbf{h}_j^{t'}} \cdot \frac{\partial \mathbf{h}_j^{t'}}{\partial W_{ji}}}_{\stackrel{\text{def}}{=} \boldsymbol{\epsilon}_{ji}^t} \quad (21)$$

608 This completes the proof of equations (1), (3), (13).

609 Derivation of eligibility traces for concrete neuron models

610 The eligibility traces for LSTMs are derived in the supplementary materials. Below we
611 provide the derivation of eligibility traces for spiking neurons.

612 **Eligibility traces for LIF neurons.** We compute the eligibility trace of a synapse of
613 a LIF neuron without adaptive threshold (equation (6)). Here the hidden state \mathbf{h}_j^t of a
614 neuron consists just of the membrane potential v_j^t and we have $\frac{\partial \mathbf{h}_j^{t+1}}{\partial \mathbf{h}_j^t} = \frac{\partial v_j^{t+1}}{\partial v_j^t} = \alpha$ and
615 $\frac{\partial v_j^t}{\partial W_{ji}} = z_i^{t-1}$ (for a derivation of the eligibility traces taking the reset into account we refer
616 to section S1.2). Using these derivatives and equation (14), one obtains that the eligibility
617 vector is the low-pass filtered presynaptic spike-train,

$$\boldsymbol{\epsilon}_{ji}^{t+1} = \mathcal{F}_\alpha(z_i^t) \stackrel{\text{def}}{=} \bar{z}_i^t, \quad (22)$$

618 and following equation (13), the eligibility trace is:

$$e_{ji}^{t+1} = \psi_j^{t+1} \bar{z}_i^t. \quad (23)$$

619 For all neurons j the derivations in the next sections also hold for synaptic connections from
620 input neurons i , but one needs to replace the network spikes z_i^{t-1} by the input spikes x_i^t (the
621 time index switches from $t-1$ to t because the hidden state $\mathbf{h}_j^t = M(\mathbf{h}_j^{t-1}, \mathbf{z}^{t-1}, \mathbf{x}^t, \mathbf{W}_j)$
622 is defined as a function of the input at time t but the preceding recurrent activity). For
623 simplicity we have focused on the case where transmission delays between neurons in the
624 RSNN are just 1 ms. If one uses more realistic length of delays d , this $-d$ appears in
625 equations (23)–(25) instead of -1 as the most relevant time point for presynaptic firing
626 (see section S1.3). This moves resulting synaptic plasticity rules closer to experimentally
627 observed forms of STDP.

628 **Eligibility traces for ALIF neurons.** The hidden state of an ALIF neuron is a two
629 dimensional vector $\mathbf{h}_j^t = [v_j^t, a_j^t]$. Hence a two dimensional eligibility vector $\boldsymbol{\epsilon}_{ji}^t \stackrel{\text{def}}{=} [\epsilon_{ji,v}^t, \epsilon_{ji,a}^t]$

630 is associated with the synapse from neuron i to neuron j , and the matrix $\frac{\partial \mathbf{h}_j^{t+1}}{\partial \mathbf{h}_j^t}$ is a 2×2
 631 matrix. The derivatives $\frac{\partial a_j^{t+1}}{\partial a_j^t}$ and $\frac{\partial a_j^{t+1}}{\partial v_j^t}$ capture the dynamics of the adaptive threshold.
 632 Hence to derive the computation of eligibility traces we substitute the spike z_j in equation
 633 (10) by its definition given in equation (9). With this convention one finds that the diagonal
 634 of the matrix $\frac{\partial \mathbf{h}_j^{t+1}}{\partial \mathbf{h}_j^t}$ is formed by the terms $\frac{\partial v_j^{t+1}}{\partial v_j^t} = \alpha$ and $\frac{\partial a_j^{t+1}}{\partial a_j^t} = \rho - \psi_j^t \beta$. Above and
 635 below the diagonal, one finds respectively $\frac{\partial v_j^{t+1}}{\partial a_j^t} = 0$, $\frac{\partial a_j^{t+1}}{\partial v_j^t} = \psi_j^t$. Seeing that $\frac{\partial \mathbf{h}_j^t}{\partial W_{ji}} =$
 636 $\left[\frac{\partial v_j^t}{\partial W_{ji}}, \frac{\partial a_j^t}{\partial W_{ji}} \right] = [z_i^{t-1}, 0]$, one can finally compute the eligibility traces using equation (13).
 637 The component of the eligibility vector associated with the membrane potential remains
 638 the same as in the LIF case and only depends on the presynaptic neuron: $\epsilon_{ji,v}^t = \bar{z}_i^{t-1}$.
 639 For the component associated with the adaptive threshold we find the following recursive
 640 update:

$$\epsilon_{ji,a}^{t+1} = \psi_j^t \bar{z}_i^{t-1} + (\rho - \psi_j^t \beta) \epsilon_{ji,a}^t, \quad (24)$$

641 and, since $\frac{\partial z_j^t}{\partial \mathbf{h}_j^t} = \left[\frac{\partial z_j^t}{\partial v_j^t}, \frac{\partial z_j^t}{\partial a_j^t} \right] = [\psi_j^t, -\beta \psi_j^t]$, this results in an eligibility trace of the form:

$$e_{ji}^t = \psi_j^t \left(\bar{z}_i^{t-1} - \beta \epsilon_{ji,a}^t \right). \quad (25)$$

642 Recall that the constant $\rho = \exp(-\frac{\delta t}{\tau_a})$ arises from the adaptation time constant τ_a , which
 643 typically lies in the range of hundreds of milliseconds to a few seconds in our experiments,
 644 yielding values of ρ between 0.995 and 0.9995. The constant β is typically of the order of
 645 0.07 in our experiments.

646 To provide a more interpretable form of eligibility trace that fits into the standard form
 647 of local terms considered in 3-factor learning rules [9], one may drop the term $-\psi_j^t \beta$ in
 648 equation (24). This approximation $\hat{\epsilon}_{ji,a}^t$ of equation (24) becomes an exponential trace
 649 of the post-pre pairings accumulated within a time window as large as the adaptation
 650 adaptation time constant:

$$\hat{\epsilon}_{ji,a}^{t+1} = \mathcal{F}_\rho \left(\psi_j^t \bar{z}_i^{t-1} \right). \quad (26)$$

651 The eligibility traces are computed with equation (24) in most experiments, but the per-
 652 formances obtained with *symmetric e-prop* and this simplification were indistinguishable
 653 in the task where temporal credit assignment is difficult of Fig. 3.

654 Synaptic plasticity rules resulting from *e-prop*

655 An exact computation of the ideal learning signal $\frac{dE}{dz_j^t}$ in equation (1) requires to back-
 656 propagate gradients through time (see Fig. 6c). For online *e-prop* we replace it with the
 657 partial derivative $\frac{\partial E}{\partial z_j^t}$, which can be computed online. Implementing the weight updates

658 with gradient descent and learning rate η , all the following plasticity rules are derived from
 659 the formula

$$\Delta W_{ji}^{\text{rec}} = -\eta \sum_t \frac{\partial E}{\partial z_j^t} e_{ji}^t. \quad (27)$$

660 Note that in the absence of the superscript t , ΔW_{ji} denotes the cumulated weight change
 661 over one trial or batch of consecutive trials but not the instantaneous weight update. This
 662 can be implemented online by accumulating weight updates in a hidden synaptic variable.
 663 Note also that the weight updates derived in the following for the recurrent weights W_{ji}^{rec}
 664 also apply to the inputs weights W_{ji}^{in} . For the output weights and biases the derivation does
 665 not require the theory of *e-prop*, and the weight updates can be found in the section S3.1.

666 **Case of regression tasks.** In the case of a regression problem with targets $y_k^{*,t}$ and out-
 667 puts y_k^t defined in equation (11), we define the loss function $E = \frac{1}{2} \sum_{t,k} (y_k^t - y_k^{*,t})^2$. This
 668 results in a partial derivative of the form $\frac{\partial E}{\partial z_j^t} = \sum_k W_{kj}^{\text{out}} \sum_{t' \geq t} (y_k^{t'} - y_k^{*,t'}) \kappa^{t'-t}$. This seem-
 669 ingly provides an obstacle for online learning, because the partial derivative is a weighted
 670 sum over future errors. But this problem can be resolved since one can interchange the
 671 two summation indices in the expression for the weight updates (see section S3.1). In this
 672 way the sum over future events transforms into a low-pass filtering of the eligibility traces
 673 $\bar{e}_{ji}^t = \mathcal{F}_\kappa(e_{ji}^t)$, and the resulting weight update can be written as

$$\Delta W_{ji}^{\text{rec}} = -\eta \sum_t \underbrace{\left(\sum_k B_{jk}(y_k^t - y_k^{*,t}) \right)}_{=L_j^t} \bar{e}_{ji}^t. \quad (28)$$

674 **Case of classification tasks.** We assume that K target categories are provided in the
 675 form of a one-hot encoded vector $\boldsymbol{\pi}^{*,t}$ with K dimensions. We define the probability for
 676 class k predicted by the network as $\pi_k^t = \text{softmax}_k(y_1^t, \dots, y_K^t) = \exp(y_k^t) / \sum_{k'} \exp(y_{k'}^t)$, and
 677 the loss function for classification tasks as the cross-entropy error $E = - \sum_{t,k} \pi_k^{*,t} \log \pi_k^t$.
 678 The plasticity rule resulting from *e-prop* reads (see derivation in section S3.1):

$$\Delta W_{ji}^{\text{rec}} = -\eta \sum_t \underbrace{\left(\sum_k B_{jk}(\pi_k^t - \pi_k^{*,t}) \right)}_{=L_j^t} \bar{e}_{ji}^t. \quad (29)$$

679 **Reward-based e-prop: application of e-prop to deep RL**

680 For reinforcement learning, the network interacts with an external environment. At any
 681 time t the environment can provide a positive or negative reward r^t . Based on the ob-
 682 servations \boldsymbol{x}^t that are perceived, the network has to commit to actions $a^{t_0}, \dots, a^{t_n}, \dots$ at
 683 certain decision times t_0, \dots, t_n, \dots . Each action a^t is sampled from a probability distri-
 684 bution $\pi(\cdot | \boldsymbol{y}^t)$ which is also referred to as the policy of the RL agent. The policy is defined

685 as function of the network outputs \mathbf{y}^t , and is chosen here to be a categorical distribution
 686 of K discrete action choices. We assume that the agent chooses action k with probability
 687 $\pi_k^t = \pi(a^t = k | \mathbf{y}^t) = \text{softmax}_k(y_1^t, \dots, y_K^t) = \exp(y_k^t) / \sum_{k'} \exp(y_{k'}^t)$.

688 The goal of reinforcement learning is to maximize the expected sum of discounted
 689 rewards. That is, we want to maximize the expected return at time $t = 0$, $\mathbb{E}[R^0]$, where
 690 the return at time t is defined as $R^t = \sum_{t' \geq t} \gamma^{t'-t} r^{t'}$ with a discount factor $\gamma \leq 1$. The
 691 expectation is taken over the agent actions a^t , the rewards r^t and the observations from the
 692 environment \mathbf{x}^t . We approach this optimization problem by using the actor-critic variant of
 693 the policy gradient algorithm, which applies gradient ascent to maximize $\mathbb{E}[R^0]$. The basis
 694 of the estimated gradient relies on an estimation of the policy gradient, as shown in section
 695 13.3 in [13]. There, the resulting weight update is given in equation (13.8), where G_t refers
 696 to the return R^t . Hence, the gradient $\frac{d\mathbb{E}[R^0]}{dW_{ji}}$ is proportional to $\mathbb{E} \left[\sum_{t_n} R^{t_n} \frac{d \log \pi(a^{t_n} | \mathbf{y}^{t_n})}{dW_{ji}} \right]$,
 697 which is easier to compute because the expectation can be estimated by an average over
 698 one or many trials. Following this strategy, we define the per-trial loss function E_π as a
 699 function of the sequence of actions $a^{t_0}, \dots, a^{t_n}, \dots$ and rewards r^0, \dots, r^T sampled during
 700 this trial:

$$E_\pi(\mathbf{z}^0, \dots, \mathbf{z}^T, a^{t_0}, \dots, a^{t_n}, \dots, r^0, \dots, r^T) \stackrel{\text{def}}{=} - \sum_n R^{t_n} \log \pi(a^{t_n} | \mathbf{y}^{t_n}) . \quad (30)$$

701 And thus:

$$\frac{d\mathbb{E}[R^0]}{dW_{ji}} \propto \mathbb{E} \left[\sum_{t_n} R^{t_n} \frac{d \log \pi(a^{t_n} | \mathbf{y}^{t_n})}{dW_{ji}} \right] = -\mathbb{E} \left[\frac{dE_\pi}{dW_{ji}} \right] . \quad (31)$$

702 Intuitively, given a trial with high rewards, policy gradient changes the network output \mathbf{y}
 703 to increase the probability of the actions a^{t_n} that occurred during this trial. In practice, the
 704 gradient $\frac{dE_\pi}{dW_{ji}}$ is known to have high variance and the efficiency of the learning algorithm
 705 can be improved using the actor-critic variant of the policy gradient algorithm. It involves
 706 the policy π (the actor) and an additional output neuron V^t which predicts the value
 707 function $\mathbb{E}[R^t]$ (the critic). The actor and the critic are learnt simultaneously by defining
 708 the loss function as

$$E = E_\pi + c_V E_V , \quad (32)$$

709 where $E_\pi = - \sum_n R^{t_n} \log \pi(a^{t_n} | \mathbf{y}^{t_n})$ measures the performance of the stochastic policy π ,
 710 and $E_V = \sum_t \frac{1}{2} (R^t - V^t)^2$ measures the accuracy of the value estimate V^t .

711 Since V^t is independent of the action a^t one can show that $0 = \mathbb{E} \left[V^{t_n} \frac{d \log \pi(a^{t_n} | \mathbf{y}^{t_n})}{dW_{ji}} \right]$.

712 We can use that to define an estimator $\widehat{\frac{dE}{dW_{ji}}}$ of the loss gradient with reduced variance:

$$-\frac{d\mathbb{E}[R^0]}{dW_{ji}} + c_V \mathbb{E} \left[\frac{dE_V}{dW_{ji}} \right] \propto \mathbb{E} \left[\frac{dE}{dW_{ji}} \right] \quad (33)$$

$$= \mathbb{E} \left[\underbrace{-\sum_{t_n} (R^{t_n} - V^{t_n}) \frac{d \log \pi(a^{t_n} | \mathbf{y}^{t_n})}{dW_{ji}} + c_V \frac{dE_V}{dW_{ji}}}_{\stackrel{\text{def}}{=} \widehat{\frac{dE}{dW_{ji}}}} \right], \quad (34)$$

713 similarly as in equation (13.11) of section 13.4 in [13]. A difference in notation is that $b(S_t)$
 714 refers to our value estimation V^t . In addition, equation (34) already includes the gradient
 715 $\frac{dE_V}{dW_{ji}}$ that is responsible for learning the value prediction. Until now this derivation follows
 716 the classical definition of the actor-critic variant of policy gradient, and the gradient $\widehat{\frac{dE}{dW_{ji}}}$
 717 can be computed with *BPTT*. To derive *reward-based e-prop* we follow instead the generic
 718 online approximation of *e-prop* as in equation (27) and approximate $\widehat{\frac{dE}{dW_{ji}}}$ by a sum of terms
 719 of the form $\widehat{\frac{\partial E}{\partial z_j^t}} e_{ji}^t$ with

$$\widehat{\frac{\partial E}{\partial z_j^t}} = -\sum_n (R^{t_n} - V^{t_n}) \frac{\partial \log \pi(a^{t_n} | \mathbf{y}^{t_n})}{\partial z_j^t} + c_V \frac{\partial E_V}{\partial z_j^t}. \quad (35)$$

720 We choose this estimator $\widehat{\frac{\partial E}{\partial z_j^t}}$ of the loss derivative because it is unbiased and has a low
 721 variance, more details are given in section S5.1. We derive below the resulting synaptic
 722 plasticity rule as needed to solve the task of Fig. 4, 5. For the case of a single action as
 723 used in Fig. S5 we refer to section S5.1.

724 When there is a delay between the action and the reward or, even harder, when a
 725 sequence of many actions lead together to a delayed reward, the loss function E cannot be
 726 computed online because the evaluation of R^{t_n} requires knowledge of future rewards. To
 727 overcome this, we introduce temporal difference errors $\delta^t = r^t + \gamma V^{t+1} - V^t$ (see Fig. 4), and
 728 use the equivalence between the forward and backward view in reinforcement learning [13].
 729 Using the one-hot encoded action $\mathbb{1}_{a^t=k}$ at time t , which assumes the value 1 if and only if
 730 $a^t = k$ (else it has value 0), we arrive at the following synaptic plasticity rules for a general
 731 actor-critic algorithm with *e-prop* (see section S5.1):

$$\Delta W_{ji}^{\text{rec}} = -\eta \sum_t \delta^t \mathcal{F}_\gamma \left(L_j^t \bar{e}_{ji}^t \right) \quad \text{for} \quad (36)$$

$$L_j^t = -c_V B_j^V + \sum_k B_{jk}^\pi (\pi_k^t - \mathbb{1}_{a^t=k}), \quad (37)$$

732 where we define the term $\pi_k^t - \mathbb{1}_{a^t=k}$ to have value zero when no action is taken at time t .
 733 B_j^V is here the weight from the output neuron for the value function to neuron j , and the
 734 weights B_{jk}^π denote the weights from the outputs for the policy.

735 A combination of reward prediction error and neuron-specific learning signal was pre-
736 viously used in a plasticity rule for feedforward networks inspired by neuroscience [58, 59].
737 Here it arises from the approximation of *BPTT* by *e-prop* in RSNNs solving RL problems.
738 Note that the filtering \mathcal{F}_γ requires an additional eligibility trace per synapse. This arises
739 from the temporal difference learning in RL [13]. It depends on the learning signal and
740 does not have the same function as the eligibility trace e_{ji}^t .

741 Code availability

742 An implementation of *e-prop* solving the tasks of Fig. 2 to 5 is made public together with
743 the publication of this paper https://github.com/IGITUGraz/eligibility_propagation.

744 Data availability

745 Data for the TIMIT and ATARI benchmark tasks were published in previous works [22, 29].
746 Data for the temporal credit assignment task are generated by a custom code provided in
747 the abovementioned code repository.

748 Acknowledgments

749 This research/project was supported by the Human Brain Project (Grand Agreement num-
750 ber 785907) and the SYNCH project (Grand Agreement number 824162) of the European
751 Union. We gratefully acknowledge the support of NVIDIA Corporation with the donation
752 of the Quadro P6000 GPU used for this research. Computations were carried out on the
753 Human Brain Project PCP Pilot Systems at the Juelich Supercomputing Centre, which
754 received co-funding from the European Union (Grand Agreement number 604102) and on
755 the Vienna Scientific Cluster (VSC).

756 We thank Thomas Bohnstingl, Wulfram Gerstner, Christopher Harvey, Martin Vinck,
757 Jason MacLean, Adam Santoro, Christopher Summerfield, and Yuqing Zhu for helpful
758 comments on an earlier version of the manuscript. Special thanks go to Arjun Rao for
759 letting us use his code for the regularization of membrane voltages.

760 **Authors contributions** GB, FS, AS and WM conceived the work, GB, FS, AS, EH
761 and DS carried out experiments and all authors contributed to the writing of the paper.

762 References

- 763 [1] LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* (2015).
- 764 [2] Allen Institute: Cell Types Database. © 2018 Allen Institute for Brain Science. Allen
765 Cell Types Database, cell feature search. Available from: celltypes.brain-map.org/data
766 (2018).

- 767 [3] Bellec, G., Salaaj, D., Subramoney, A., Legenstein, R. & Maass, W. Long short-term
768 memory and learning-to-learn in networks of spiking neurons. *NeurIPS* (2018).
- 769 [4] Huh, D. & Sejnowski, T. J. Gradient descent for spiking neural networks. *NeurIPS*
770 (2018).
- 771 [5] Lillicrap, T. P. & Santoro, A. Backpropagation through time and the brain. *Current*
772 *Opinion in Neurobiology* (2019).
- 773 [6] Sanhueza, M. & Lisman, J. The CAMKII/NMDAR complex as a molecular memory.
774 *Molecular Brain* **6**, 10 (2013).
- 775 [7] Cassenaer, S. & Laurent, G. Conditional modulation of spike-timing-dependent plas-
776 ticity for olfactory learning. *Nature* (2012).
- 777 [8] Yagishita, S. *et al.* A critical time window for dopamine actions on the structural
778 plasticity of dendritic spines. *Science* (2014).
- 779 [9] Gerstner, W., Lehmann, M., Liakoni, V., Corneil, D. & Brea, J. Eligibility Traces and
780 Plasticity on Behavioral Time Scales: Experimental Support of NeoHebbian Three-
781 Factor Learning Rules. *Frontiers in Neural Circuits* (2018).
- 782 [10] Sajad, A., Godlove, D. C. & Schall, J. D. Cortical microcircuitry of performance
783 monitoring. *Nature Neuroscience* (2019).
- 784 [11] Engelhard, B. *et al.* Specialized coding of sensory, motor and cognitive variables in
785 VTA dopamine neurons. *Nature* (2019).
- 786 [12] Roeper, J. Dissecting the diversity of midbrain dopamine neurons. *Trends in neuro-*
787 *sciences* (2013).
- 788 [13] Sutton, R. S. & Barto, A. G. *Reinforcement Learning: An Introduction* (MIT press,
789 2018).
- 790 [14] Mnih, V. *et al.* Asynchronous methods for deep reinforcement learning. In *ICML*,
791 1928–1937 (2016).
- 792 [15] Williams, R. J. & Zipser, D. A learning algorithm for continually running fully recur-
793 rent neural networks. *Neural computation* **1**, 270–280 (1989).
- 794 [16] Furber, S. B., Galluppi, F., Temple, S. & Plana, L. A. The SpiNNaker project.
795 *Proceedings of the IEEE* **102**, 652–665 (2014).
- 796 [17] Davies, M. *et al.* Loihi: A neuromorphic manycore processor with on-chip learning.
797 *IEEE Micro* (2018).

- 798 [18] Lillicrap, T. P., Cownden, D., Tweed, D. B. & Akerman, C. J. Random synaptic
799 feedback weights support error backpropagation for deep learning. *Nature Communi-*
800 *cations* (2016).
- 801 [19] Nøkland, A. Direct feedback alignment provides learning in deep neural networks. In
802 *NIPS* (2016).
- 803 [20] Samadi, A., Lillicrap, T. P. & Tweed, D. B. Deep learning with dynamic spiking
804 neurons and fixed feedback weights. *Neural computation* **29**, 578–602 (2017).
- 805 [21] Clopath, C., Büsing, L., Vasilaki, E. & Gerstner, W. Connectivity reflects coding: a
806 model of voltage-based STDP with homeostasis. *Nature Neuroscience* (2010).
- 807 [22] Garofolo, J. S., Lamel, L. F., Fisher, W. M., Fiscus, J. G. & Pallett, D. S. DARPA
808 TIMIT acoustic-phonetic continuous speech corpus CD-ROM. *NASA STI/Recon Tech-*
809 *nical Report N* (1993).
- 810 [23] Greff, K., Srivastava, R. K., Koutník, J., Steunebrink, B. R. & Schmidhuber, J. LSTM:
811 A search space odyssey. *IEEE TNNLS* (2017).
- 812 [24] Graves, A., Mohamed, A.-R. & Hinton, G. Speech recognition with deep recurrent
813 neural networks. *ICASSP* (2013).
- 814 [25] Morcos, A. S. & Harvey, C. D. History-dependent variability in population dynamics
815 during evidence accumulation in cortex. *Nature Neuroscience* (2016).
- 816 [26] Kappel, D., Legenstein, R., Habenschuss, S., Hsieh, M. & Maass, W. A dynamic
817 connectome supports the emergence of stable computational function of neural circuits
818 through reward-based learning. *eNeuro* (2018).
- 819 [27] Bartunov, S. *et al.* Assessing the scalability of biologically-motivated deep learning
820 algorithms and architectures. In *Advances in Neural Information Processing Systems*
821 (2018).
- 822 [28] Frémaux, N. & Gerstner, W. Neuromodulated spike-timing-dependent plasticity, and
823 theory of three-factor learning rules. *Frontiers in neural circuits* **9**, 85 (2016).
- 824 [29] Bellemare, M. G., Naddaf, Y., Veness, J. & Bowling, M. The arcade learning envi-
825 ronment: An evaluation platform for general agents. *Journal of Artificial Intelligence*
826 *Research* **47**, 253–279 (2013).
- 827 [30] Runyan, C. A., Piasini, E., Panzeri, S. & Harvey, C. D. Distinct timescales of popu-
828 lation coding across cortex. *Nature* (2017).
- 829 [31] Bellec, G. *et al.* Biologically inspired alternatives to backpropagation through time
830 for learning in recurrent neural nets. *arXiv:1901.09049* (2019).

- 831 [32] Gilra, A. & Gerstner, W. Predicting non-linear dynamics by stable local learning in
832 a recurrent spiking neural network. *Elife* **6**, e28295 (2017).
- 833 [33] Thalmeier, D., Uhlmann, M., Kappen, H. J. & Memmesheimer, R.-M. Learning
834 universal computations with spikes. *PLoS computational biology* **12** (2016).
- 835 [34] Alemi, A., Machens, C. K., Deneve, S. & Slotine, J.-J. Learning nonlinear dynamics
836 in efficient, balanced spiking networks using local plasticity rules. In *Thirty-Second
837 AAAI Conference on Artificial Intelligence* (2018).
- 838 [35] Nicola, W. & Clopath, C. Supervised learning in spiking neural networks with force
839 training. *Nature Communications* (2017).
- 840 [36] Sussillo, D. & Abbott, L. F. Generating coherent patterns of activity from chaotic
841 neural networks. *Neuron* **63**, 544–557 (2009).
- 842 [37] Abbott, L. F., DePasquale, B. & Memmesheimer, R.-M. Building functional networks
843 of spiking model neurons. *Nature neuroscience* **19**, 350 (2016).
- 844 [38] Ingrosso, A. & Abbott, L. Training dynamically balanced excitatory-inhibitory net-
845 works. *PloS one* **14** (2019).
- 846 [39] Kim, C. M. & Chow, C. C. Learning recurrent dynamics in spiking networks. *eLife*
847 **7**, e37124 (2018).
- 848 [40] Zenke, F. & Ganguli, S. Superspike: Supervised learning in multilayer spiking neural
849 networks. *Neural computation* (2018).
- 850 [41] Shrestha, S. B. & Orchard, G. Slayer: Spike layer error reassignment in time. In
851 Bengio, S. *et al.* (eds.) *NeurIPS* (2018).
- 852 [42] Nefteci, E. O., Augustine, C., Paul, S. & Detorakis, G. Event-driven random back-
853 propagation: Enabling neuromorphic deep learning machines. *Frontiers in neuro-
854 science* **11**, 324 (2017).
- 855 [43] Kaiser, J., Mostafa, H. & Nefteci, E. Synaptic plasticity dynamics for deep continuous
856 local learning. *arXiv preprint arXiv:1811.10766* (2018).
- 857 [44] Emre O. Nefteci, F. Z., Hesham Mostafa. Surrogate gradient learning in spiking neu-
858 ral networks: Bringing the power of gradient-based optimization to spiking neural
859 networks. *IEEE Signal Processing Magazine* (2019).
- 860 [45] Bengio, Y., Léonard, N. & Courville, A. Estimating or propagating gradients through
861 stochastic neurons for conditional computation. *arXiv preprint arXiv:1308.3432*
862 (2013).

- 863 [46] Courbariaux, M., Hubara, I., Soudry, D., El-Yaniv, R. & Bengio, Y. Binarized neural
864 networks: Training deep neural networks with weights and activations constrained
865 to+ 1 or-1. *arXiv preprint arXiv:1602.02830* (2016).
- 866 [47] Marschall, O., Cho, K. & Savin, C. A unified framework of online learning algorithms
867 for training recurrent neural networks. *arXiv preprint arXiv:1907.02649* (2019).
- 868 [48] Mujika, A., Meier, F. & Steger, A. Approximating real-time recurrent learning with
869 random kronecker factors. *NeurIPS* (2018).
- 870 [49] Tallec, C. & Ollivier, Y. Unbiased online recurrent optimization. *ICLR* (2018).
- 871 [50] Roth, C., Kanitscheider, I. & Fiete, I. Kernel rnn learning (kernel). *ICLR* (2019).
- 872 [51] Murray, J. M. Local online learning in recurrent networks with random feedback.
873 *eLife* (2019).
- 874 [52] Jaderberg, M. *et al.* Decoupled neural interfaces using synthetic gradients. *arXiv*
875 *preprint arXiv:1608.05343* (2016).
- 876 [53] Gerstner, W., Kistler, W. M., Naud, R. & Paninski, L. *Neuronal dynamics: From sin-*
877 *gle neurons to networks and models of cognition* (Cambridge University Press, 2014).
- 878 [54] Pozzorini, C. *et al.* Automated high-throughput characterization of single neurons by
879 means of simplified spiking models. *PLoS Computational Biology* (2015).
- 880 [55] Gouwens, N. W. *et al.* Systematic generation of biophysically detailed models for
881 diverse cortical neuron types. *Nature Communications* (2018).
- 882 [56] Esser, S. K. *et al.* Convolutional networks for fast, energy-efficient neuromorphic
883 computing. *PNAS* (2016).
- 884 [57] Werbos, P. J. Backpropagation through time: what it does and how to do it. *Pro-*
885 *ceedings of the IEEE* (1990).
- 886 [58] Roelfsema, P. R. & Holtmaat, A. Control of synaptic plasticity in deep cortical
887 networks. *Nature Reviews Neuroscience* (2018).
- 888 [59] Pozzi, I., Bohté, S. & Roelfsema, P. A biologically plausible learning rule for deep
889 learning in the brain. *arXiv preprint arXiv:1811.01768* (2018).