

Deep reinforcement learning in a time-continuous model

Akos F. Kungl¹, Dominik Dold¹, Oskar Riedler², Walter Senn³,
Mihai A Petrovici^{1,3}

¹) Heidelberg University, Kirchhoff Institute for Physics

²) Heidelberg University, Faculty of Mathematics and Computer Science

³) University of Bern, Department of Physiology

Inspired by the recent success of deep learning [1], several models emerged trying to explain how the brain might realize plasticity rules reaching similar performances as deep learning [2, 3, 4, 5]. However, all of these models consider only supervised and unsupervised learning, where an external teacher is needed to produce an error signal that guides plasticity.

In this work, we introduce a model of reinforcement learning with the principle of Neuronal Least Action (R-NLA). We extend previous works on time-continuous error backpropagation in cortical microcircuits [4, 6] to achieve a biologically plausible model implementing deep reinforcement learning.

In R-NLA the neurosynaptic dynamics is derived from the energy function using the variational principle. In the resulting dynamics the phase-advanced firing of the neurons effectively undoes the network delay introduced by finite membrane time-constants. Errors are introduced to the network by nudging, and they are propagated to deeper layers via cortical microcircuits. Instead of having an explicit teacher, the output neurons, which represent the actions, form a soft winner-take-all network (Fig A). This winner-take-all structure evokes a nudging on the soma of the output neurons, which is subsequently backpropagated through the network. A reward prediction error $\delta = R - \langle R \rangle$ modulates the plasticity multiplicatively as a formally deduced global reward-specific neuromodulator [7]. By construction, the learning rule approximates the policy gradient of the mean expected reward.

We show, on a toy problem, that R-NLA can learn classification tasks in the reinforcement learning framework with similar performance as an equivalent deep reinforcement learning model (Fig B). Further, we show that it is robust against time delayed rewards, even if the reward-delay is not constant but randomly distributed (Fig C).

R-NLA constitutes a time-continuous implementation of biologically plausible deep reinforcement learning, robust to delayed reward. The self-teaching soft winner-take-all mechanism removes the necessity of an explicit teacher and the proposed learning rule solves the problem of synaptic consolidation. The model can be extended to an actor-critic model, where a second (deep) critic network learns the state-value function.

References

- [1] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436, 2015.

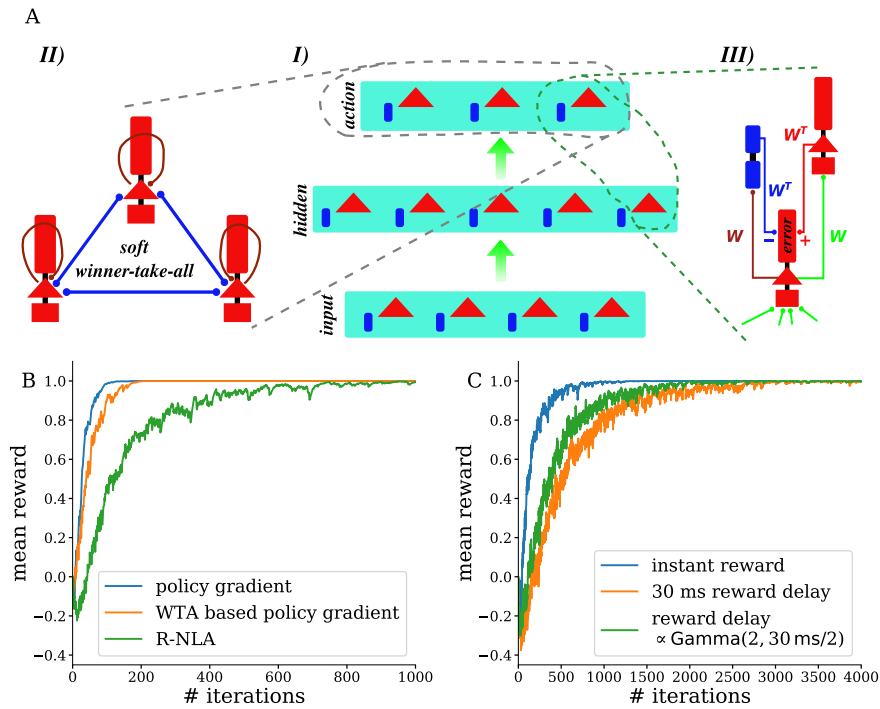


Figure 1: **A-I** Network schematics. **A-II** Soft winner-take-all network in the output layer. **A-III** Microcircuit for error backpropagation. **B** Comparison to classical reinforcement learning methods. **C** Robustness with respect to reward delays.

- [2] Randall C O'Reilly. Biologically plausible error-driven learning using local activation differences: The generalized recirculation algorithm. *Neural Computation*, 8(5):895–938, 1996.
- [3] James CR Whittington and Rafal Bogacz. An approximation of the error backpropagation algorithm in a predictive coding network with local hebbian synaptic plasticity. *Neural Computation*, 29(5):1229–1262, 2017.
- [4] João Sacramento, Rui Ponte Costa, Yoshua Bengio, and Walter Senn. Dendritic cortical microcircuits approximate the backpropagation algorithm. In *Advances in Neural Information Processing Systems*, pages 8721–8732, 2018.
- [5] James CR Whittington and Rafal Bogacz. Theories of error backpropagation in the brain. *Trends in Cognitive Sciences*, 2019.
- [6] Dominik Dold, Akos F Kungl, João Sacramento, Mihai A Petrovici, Kaspar Schindler, Jonathan Binas, Yoshua Bengio, and Walter Senn. Lagrangian dynamics of dendritic microcircuits enables real-time backpropagation of errors. In *Cosyne Conference*, 2019.

- [7] Nicolas Frémaux and Wulfram Gerstner. Neuromodulated spike-timing-dependent plasticity, and theory of three-factor learning rules. *Frontiers in Neural Circuits*, 9:85, 2016.