

1 We thank the reviewers for their constructive feedback.

2 Reviewer 2 pointed out that cross frequency coupling is “speculative” but this topic has been discussed widely in
3 particular in Parkinson which is cited also in this paper. There is also a very recent modelling paper:

4 Belić JJ, Halje P, Richter U, Petersson P and Hellgren Kotaleski J (2016) Untangling Cortico-Striatal
5 Connectivity and Cross-Frequency Coupling in L-DOPA-Induced Dyskinesia. Front. Syst. Neurosci.
6 10:26. doi: 10.3389/fnsys.2016.00026

7 However, we agree with the reviewer that more neurophysiological work needs to be done to explore cross frequency
8 coupling and plasticity. Clearly the brain uses both high frequency bands which can easily change plasticity and low
9 frequency bands which cannot.

10 Reviewer 1/3: H_0 is simply the transfer function of the fixed feedback loop – e.g. if it were a simple thermostat trying
11 to maintain a desired temperature, it would be a threshold function that mapped temperatures onto a binary output. P_1
12 is the function that transforms the actions back into sensory inputs V_i used for learning. Happy to add this to the text.

13 Reviewers 1/3: we have confused you with the unexplained term ‘reflex’, apologies! In figure 1, we see that there is
14 a fixed feedback loop with setpoint SP; this is what we are calling a ‘reflex’, as an analogy with the reflex circuits
15 that are often seen in biology. The reflex for the shooter game is rather artificial, but reflexes are common biological
16 control mechanisms for coping with disturbances. However they have the drawback that they are purely reactive. The
17 contribution here is:

- 18 • To show that the reflex can also provide a learning signal for training a deep network, and that network is then
19 learning input control rather than output control - it is learning to keep the reflex silent rather than produce
20 target outputs.
- 21 • To achieve learning where the errors are defined at the inputs, we develop a learning rule is not gradient-based,
22 but rather where both activity and error signals are propagated forwards in the same weighted fashion. The
23 learning rule is then correlation-based between these two signals.
- 24 • The “reflex” was chosen because it’s possible to treat it analytically with control theory. However, the error
25 can also originate from, for example the “critic” of TD learning aka the “reward prediction error” (and is thus
26 also biologically realistic). This would indeed allow more predictive tasks such as Atari, etc but is far beyond
27 the scope to show a very new concept as pointed out by all reviewers.

28 All reviewers, in particular reviewer 1 about Figure 2B: We thought that diagram would help to understand the learning
29 Eq. 2&3, where we correlate the activity v_j at neuron j with the error signal e_k at neuron k . This is identical to
30 backpropagation of course - the key difference lies in the fact that the error is calculated at the inputs and propagated
31 forwards in a weighted fashion (Eq.3), as is the case in closed loop systems. We’d be very happy to expand this section
32 to improve the clarity of the algorithm as suggested by reviewer 1 and tie this in with error feedback.

33 Reviewer 1 about behavioural flexibility: The reason we would expect more behavioural flexibility from this approach
34 is that learning does not explicitly evaluate the network outputs at all. Of course, the inputs do relate to the outputs, but
35 only indirectly, and only via the transfer function $P_{0/1}$ of the environment. So the outputs are less constrained than they
36 would be under Q-learning for example, where the Q value is a direct function of the outputs (and the world state). We
37 should probably say “in principle this leads to greater flexibility”, until we show evidence of this.

38 Reviewer 3: The goal was to show the feasibility of an algorithm that can function in an arbitrarily deep network. It is
39 true that the actual network used was not particularly deep however. In section 4, we say “let us assume two layers j
40 and k ”, but this is merely to illustrate how the learning in one layer relates to its neighbouring layer to demonstrate a
41 novel concept.

42 Reviewer 3: For the shooter example, unfortunately some details have been lost during editing. There are 3 output
43 neurons with \tanh activation; a negative output means ‘rotate left’, a positive ‘rotate right’. The 3 neurons operate
44 with different gains, to give finer control. The terms g_{net} and g_{err} are the gains for the reflex, and the learned control
45 signal. Momentum was 0.5; weights initialised in a zero-mean uniform distribution, scaled by the number of weights
46 outputting each neuron.

47 Reviewer 1, Equation 4: The term ρ is the gain of the fixed feedback loop (the reflex), and can be absorbed into H_0 .

48 Reviewer 1, Equation 5:

$$\Delta w_{ij} = X_0(z)V_i(-z) \quad (1)$$

49 is simply a correlation of X_0 with V_i , and . The term $(-z)$ appears because correlation involves reversing one of the
50 signals in time.