

Supplementary Material S1: Conditions for a Constant Dopaminergic Baseline Firing Rate

It has been found that midbrain dopamine neurons react to a reward or reward predicting stimulus. As long as nothing unpredictable happens, the neurons fire with a constant baseline activity [1]. To incorporate this finding into our model, we need to find conditions such that the dopaminergic neurons fire with a constant baseline activity, independent of the value of the agent's current state, except during phasic activation following a state transition. Here, we derive two relationships that enable the network to fulfill this requirement: firstly between the number of neurons in the ventral pallidum N_{VP} and the number of neurons in the striatum N_{STR} , secondly between the weights of the synapses that connect the striatum to the dopaminergic neurons w_{DA}^{STR} and those that connect the ventral pallidum to the dopaminergic neurons w_{DA}^{VP} .

The firing rate λ of an integrate-and-fire neuron can be calculated from the first passage time of the membrane potential across the threshold (see [2] for a review). For leaky integrate-and-fire neuron models, provided the distribution of the free membrane potential is sufficiently close to a Gaussian, the mean μ and the variance σ^2 of this distribution determine the first passage time and therefore also the firing rate:

$$\lambda = f(\mu, \sigma^2)$$

If an agent is stationary in a state with no associated external reward, each dopamine neuron integrates inputs from three sources: N_{STR} striatal neurons, each firing at λ_{STR} with mean synaptic weight w_{DA}^{STR} , N_{VP} neurons of the ventral pallidum, each firing at λ_{VP} with mean synaptic weight w_{DA}^{VP} , and the background noise. This last term is made up of an excitatory and an inhibitory Poissonian input with firing rates $\lambda_{bg,E} = \lambda_{bg,I} = \lambda_{bg}$ and synaptic strengths $w_{bg,E} = -w_{bg,I}$. Under the assumption that all input spike trains are uncorrelated, stationary Poisson processes, the mean membrane potential μ and the variance σ^2 can be calculated according to shot noise theory [3]:

$$\begin{aligned} \mu &= N_{STR}\lambda_{STR} \int_{-\infty}^{\infty} h(w_{DA}^{STR}, t) dt + N_{VP}\lambda_{VP} \int_{-\infty}^{\infty} h(w_{DA}^{VP}, t) dt \\ \sigma^2 &= N_{STR}\lambda_{STR} \int_{-\infty}^{\infty} h^2(w_{DA}^{STR}, t) dt + N_{VP}\lambda_{VP} \int_{-\infty}^{\infty} h^2(w_{DA}^{VP}, t) dt + \lambda_{bg,E} \int_{-\infty}^{\infty} h^2(w_{bg,E}, t) dt + \lambda_{bg,I} \int_{-\infty}^{\infty} h^2(w_{bg,I}, t) dt, \end{aligned}$$

where $h(w, t)$ is the post-synaptic potential. For α -shaped postsynaptic currents we have:

$$h(w, t) = \begin{cases} \frac{ew}{\tau_\alpha C} \left(\frac{1}{\tau_m} - \frac{1}{\tau_\alpha} \right)^{-2} \left(\left[\frac{1}{\tau_m} - \frac{1}{\tau_\alpha} \right] t e^{-t/\tau_\alpha} - e^{-t/\tau_\alpha} + e^{-t/\tau_m} \right) & t \geq 0 \\ 0 & t < 0 \end{cases}$$

Therefore, the mean membrane potential and variance of a dopaminergic neuron are given by:

$$\begin{aligned} \mu &= (N_{\text{STR}} \lambda_{\text{STR}} w_{\text{DA}}^{\text{STR}} + N_{\text{VP}} \lambda_{\text{VP}} w_{\text{DA}}^{\text{VP}}) F_1 \\ \sigma^2 &= \left(N_{\text{STR}} \lambda_{\text{STR}} (w_{\text{DA}}^{\text{STR}})^2 + N_{\text{VP}} \lambda_{\text{VP}} (w_{\text{DA}}^{\text{VP}})^2 + 2\lambda_{\text{bg}} w_{\text{bg}}^2 \right) F_2 \end{aligned}$$

with

$$\begin{aligned} F_1 &= \frac{1}{C} e^{\tau_m \tau_\alpha} \\ F_2 &= \frac{1}{C^2} e^2 \tau_m^2 \tau_\alpha^2 \frac{(\tau_\alpha + 2\tau_m)}{4(\tau_\alpha + \tau_m)^2}. \end{aligned}$$

We assume a linear relationship between the firing rates of the neurons in the ventral pallidum and striatal populations:

$$\lambda_{\text{VP}} = -a\lambda_{\text{STR}} + b$$

Therefore the mean membrane potential and variance of a dopaminergic neuron whilst the agent remains in an unrewarded state are given by:

$$\begin{aligned} \mu &= \left((w_{\text{DA}}^{\text{STR}} N_{\text{STR}} - a w_{\text{DA}}^{\text{VP}} N_{\text{VP}}) \lambda_{\text{STR}} + b w_{\text{DA}}^{\text{VP}} N_{\text{VP}} \right) F_1 \\ \sigma^2 &= \left(\left((w_{\text{DA}}^{\text{STR}})^2 N_{\text{STR}} - a (w_{\text{DA}}^{\text{VP}})^2 N_{\text{VP}} \right) \lambda_{\text{STR}} + b (w_{\text{DA}}^{\text{VP}})^2 N_{\text{VP}} + 2\lambda_{\text{bg}} w_{\text{bg}}^2 \right) F_2 \end{aligned}$$

This results in the following conditions under which the dopaminergic firing rate is independent of the striatal firing rate λ_{STR} and hence of the value of the agent's current state:

$$\begin{aligned} w_{\text{DA}}^{\text{VP}} &= \frac{1}{a} w_{\text{DA}}^{\text{STR}} \frac{N_{\text{STR}}}{N_{\text{VP}}} \\ (w_{\text{DA}}^{\text{VP}})^2 &= \frac{1}{a} (w_{\text{DA}}^{\text{STR}})^2 \frac{N_{\text{STR}}}{N_{\text{VP}}} \end{aligned}$$

These conditions are fulfilled for

$$N_{\text{STR}} = aN_{\text{VP}}$$

and

$$w_{\text{DA}}^{\text{VP}} = w_{\text{DA}}^{\text{STR}}.$$

For the parameters used in our simulations, we have verified that the linear relationship between the firing rates of the neurons in the ventral pallidum and striatal populations holds with $a = 1$ resulting in a constant dopaminergic rate.

In such a network, when the agent moves from one state s_i to another state s_{i+1} the mean membrane potential and variance are given by:

$$\begin{aligned}\mu &= (-w_{\text{DA}}^{\text{STR}} N_{\text{STR}} [\lambda_{\text{STR}}(s_{i+1}) - \lambda_{\text{STR}}(s_i)] + bw_{\text{DA}}^{\text{VP}} N_{\text{VP}}) F_1 \\ \sigma^2 &= \left(- (w_{\text{DA}}^{\text{STR}})^2 N_{\text{STR}} [\lambda_{\text{STR}}(s_{i+1}) - \lambda_{\text{STR}}(s_i)] + b (w_{\text{DA}}^{\text{VP}})^2 N_{\text{VP}} + 2\lambda_{\text{bg}} w_{\text{bg}}^2 \right) F_2,\end{aligned}$$

where $\lambda_{\text{STR}}(s_x)$ is the striatal firing rate when the agent is in state s_x . Thus, the mean membrane potential as well as the variance depend on the undiscounted difference between the values of two successive states as encoded by the striatal firing rate. If the network is not tuned to result in a constant tonic dopaminergic rate, the mean membrane potential and variance can be written in the following way:

$$\begin{aligned}\mu &= \left(- [\gamma \lambda_{\text{STR}}(s_{i+1}) - \lambda_{\text{STR}}(s_i)] + \frac{bw_{\text{DA}}^{\text{VP}} N_{\text{VP}}}{w_{\text{DA}}^{\text{STR}} N_{\text{STR}}} \right) F_1 \\ \sigma^2 &= \left(- \left[\frac{w_{\text{DA}}^{\text{VP}}}{w_{\text{DA}}^{\text{STR}}} \gamma \lambda_{\text{STR}}(s_{i+1}) - \lambda_{\text{STR}}(s_i) \right] + \frac{b (w_{\text{DA}}^{\text{VP}})^2 N_{\text{VP}} + 2\lambda_{\text{bg}} w_{\text{bg}}^2}{(w_{\text{DA}}^{\text{STR}})^2 N_{\text{STR}}} \right) F_2\end{aligned}$$

with $\gamma = (aw_{\text{DA}}^{\text{VP}} N_{\text{VP}}) / (w_{\text{DA}}^{\text{STR}} N_{\text{STR}})$, i.e. they depend on the discounted difference between two states. For $0 \leq \gamma \leq 1$ the phasic dopaminergic signal could therefore be used to drive temporal-difference learning in the cortico-striatal synapses on the basis of a simplified synaptic plasticity dynamics that does not compensate for a missing γ -factor (eq. 7 in Sec. 2.3 of the main text).

References

1. Ljungberg T, Apicella P, Schultz W (1992) Responses of monkey dopamine neurons during learning of behavioral reactions. *J Neurophysiol* 67: 145–163.
2. Burkitt AN (2006) A review on the integrate-and-fire neuron model: I. homogenous synaptic input. *bicy* 95: 1–19.

3. Papoulis A (1991) Probability, Random Variables, and Stochastic Processes. Boston, Massachusetts: McGraw-Hill, 3 edition.