

1 Thanks for the comments. There were some suggested improvements – if accepted we will incorporate: an algorithm  
2 box; analysis of complexity; redesign and/or annotation for Figure 1; Figure 2A adaptation; and typo corrections.

3 **R1: Fig 2:** We will delete the reference to a gray line, which had been removed from the figure – sorry. Parameters  
4  $t_0$  and  $\tau$  affect readout, not the internal model or inference:  $t_0$  is the perceived time of the flash;  $\tau$  is the postdictive  
5 window – the duration of subsequent evidence that contributes to the estimate of the location of the moving object at  $t_0$ .

6 **Comparisons:** Biological plausibility is at the heart of the problem we seek to address. Postdictive behaviour depends  
7 on updating beliefs about *past* states – the central question of our work is how this might happen in natural online  
8 inference. Standard temporal variational methods (e.g. Chung et al., 2015; Fraccaro et al., 2017) represent states  
9 individually, and so would require acausal (“backward”) message passing to revise past states (see also comments  
10 under Signatures below). This is possible (e.g. Fraccaro; or Johnson et al., 2016), but currently only implemented for  
11 restricted classes of generative model which do not capture the natural environment.

12 **R2: Encoding functions:** indeed, we chose fixed but random  $\gamma$ . Results are robust to redrawing random projections,  
13 and to different nonlinearities. To train the parameters of  $\gamma$  and  $\psi_t$  one could minimize the predictive error (5) or (9)  
14 wrt these parameters. Gradient descent would be a natural choice, albeit with debatable biological plausibility.

15 **Signatures of internal representation:** a general attempt to distinguish encoding schemes that use expectations (DDC),  
16 samples, or natural parameters (including the ‘PPC’) depends on additional model assumptions and is larger than can  
17 be addressed satisfactorily in a NeurIPS paper. However, in the specific context of temporal models and postdiction,  
18 these schemes differ in the mechanism by which beliefs about the past may be maintained and updated. The linearity of  
19 expectation, combined with temporal encoding functions (or sufficient statistics)  $\psi_t$ , allows a simple dynamical update  
20 rule to maintain and update beliefs (eqs (7) and (10) or (11)). If beliefs about latent states were represented by samples,  
21 then these samples would need to be maintained in a sort of “buffer” to be modified by later inputs and accessed by  
22 downstream processing. Natural parameter encodings could, in principle, exploit temporal sufficient statistic functions  
23 similar to those we introduce, but the dynamical updates required to maintain and revise the natural parameters through  
24 time would be far more complex to derive and implement than for the mean parameters of the DDC. Thus, of these  
25 options, we believe the DDC inference and postdiction in the temporal context comes closest to the known biology.

26 **Lines 153-154:** (We will clarify the entire first paragraph of 3.2 if accepted.) Equation (8) would have been better  
27 written:  $\mathcal{L}^f(\widetilde{\mathbf{W}}_t; \mathbf{x}_{1:t-1}) = \mathbb{E}_{q(\mathbf{z}_{1:t}, \mathbf{x}_t | \mathbf{x}_{1:t-1})} [\|\widetilde{\mathbf{W}}_t \boldsymbol{\sigma}(\mathbf{x}_t) - \boldsymbol{\psi}(\mathbf{z}_{1:t})\|_2^2]$ . That is, the loss is history-dependent and the  
28 optimal weights thus depend on both history and time. Now, if we restrict functions  $h_{\mathbf{W}}$  in (9) to be linear in  $\boldsymbol{\sigma}(\mathbf{x}_t)$  (as  
29 both (10) and (11) are) then  $h_{\mathbf{W}}$  can be seen as effectively mapping  $\mathbf{r}_{t-1}$  to a time- and history-dependent  $\widetilde{\mathbf{W}}_t$ . Thus,  
30 although the loss (9) is the expectation of (8) over histories, the optimal sequence of  $\widetilde{\mathbf{W}}_t$  for each history still provides a  
31 target for the minimisation of (9), and the expectation of the minima of (8) bounds the minimum of (9) below. But  
32 indeed, the constrained form of  $h_{\mathbf{W}}(\mathbf{r}_{t-1}, \cdot)$  means that this lower bound will not generally be achieved.

33 **R3: “Neurons”:** Our intention is for activities of neurons in our scheme to provide a model for the firing rates of  
34 biological neurons. We expect DDC-based computations to be robust to Poisson-like noise, as linear operations acting  
35 on populations of neurons (inference and readout) effectively average away independent perturbations. Consistent with  
36 this intuition, when we re-ran the flash-lag experiment using Poisson neurons the mean results were essentially the  
37 same, albeit with greater trial-to-trial variance. For occluded tracking, the  $R^2$  in postdicting the true  $\mathbf{z}_t$  by posterior  
38 mean still increased with postdictive window  $\tau$ , although this improvement did not extend as far into the past as in the  
39 noiseless case. If accepted, we will include these new experimental results.

40 **Line 92:** Indeed, DDC activities encode a *distribution* or *belief* about a random variable – we will amend the text.

41 **Line 118:** Our focus here is to propose how neural circuits may implement postdictive inference – itself, a non-trivial  
42 problem even when the internal model of the world is known – and relate such inference to psychometric phenomena.  
43 A general approach to modelling learning using DDC-based inference in wake-sleep has been proposed before (Vertes  
44 and Sahani, 2018), although extensions to temporal models are the subject of current research.

45 **Learning the readout:** the key feature of (3) and (12) is that  $\alpha$  can be found for known target functions using  
46 evaluations of [e.g. for (12)]  $\boldsymbol{\psi}(\mathbf{z}_{1:t})$  and  $l(\mathbf{z}_{t-\tau})$ , and distributional uncertainty then propagates appropriately. Thus (as  
47 with the recognition weights) we need only simulations of  $\mathbf{z}_{1:t}$  from the internal model, potentially from the sleep phase  
48 of wake-sleep. More generally, the target function  $l(\cdot)$  might need to be learnt based on supervision or reinforcement.  
49 This could follow “normal” learning rules; with the DDC again ensuring that uncertainty is handled correctly.

50 **Lines 190-195:** These simulations probe the quality of inference in the learnt model. Figure 1D reproduces the known  
51 effect that perceptual continuity depends on the level of the interrupting noise. Figure 1E demonstrates quantitative  
52 inference – the posterior on continuity is reduced when the noise is quieter still. Finally, 1F demonstrates that the  
53 inference model has correctly encoded the fact that tones in the generative model do not change in loudness.

54 **Fig 1, multiple buses:** The hallmark of postdiction is that beliefs about past values are revised by new evidence. Thus,  
55 we show beliefs at time step  $i$  (column) conditioned on sensory experience up to various later times  $j$  (row). That is, the  
56 three boxes in the  $i$ th column of the  $j$ th row show decoded posterior probabilities for the three possible levels of the tone  
57 at step  $j$ , based on sensory input up to time  $i$ . The smaller inset buses show the inferred noise level in the same way.