

# A neurally plausible model for online recognition and postdiction: Supplementary material

## A Properties of EMSE minimization

### A.1 EMSE and the expected error in estimating the DDC of exact posterior

The following proposition sets up the relationship between the EMSE in (5) and the error between  $\mathbf{h}(\mathbf{x})$  and the true posterior mean (3).

**Proposition 1.** *Minimizing (5) minimizes  $\mathbb{E}_{p(\mathbf{x})}[\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\|\boldsymbol{\gamma}(\mathbf{z}) - \mathbf{h}(\mathbf{x})\|^2]]$ , the expected l-2 squared distance between the true posterior DDC and the recognition model prediction. The minimum is achieved when  $\mathbf{h}(\mathbf{x}) = \mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})]$*

*Proof.* Following the standard decomposition of the MSE for regression with additional expectation on  $p(\mathbf{x})$ ,

$$\begin{aligned} (5) &= \mathbb{E}_{p(\mathbf{z}, \mathbf{x})}[\|\boldsymbol{\gamma}(\mathbf{z}) - \mathbf{h}(\mathbf{x})\|^2] \\ &= \mathbb{E}_{p(\mathbf{z}, \mathbf{x})}[\|\boldsymbol{\gamma}(\mathbf{z}) - \mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})] + \mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})] - \mathbf{h}(\mathbf{x})\|^2] \\ &= \mathbb{E}_{p(\mathbf{x})} \text{Tr}[\mathbb{C}_{p(\mathbf{z}|\mathbf{x})}(\boldsymbol{\gamma}(\mathbf{z}))] + \mathbb{E}_{p(\mathbf{x})}[\|\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})] - \mathbf{h}(\mathbf{x})\|^2]. \end{aligned} \quad (15)$$

The cross term in the second line is zero because

$$\begin{aligned} &\mathbb{E}_{p(\mathbf{z}, \mathbf{x})}[(\boldsymbol{\gamma}(\mathbf{z}) - \mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})]) \cdot (\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})] - \mathbf{h}(\mathbf{x}))] \\ &= \mathbb{E}_{p(\mathbf{x})}[\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[(\boldsymbol{\gamma}(\mathbf{z}) - \mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})]) \cdot (\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})] - \mathbf{h}(\mathbf{x}))]] = 0 \end{aligned}$$

The first term is a positive constant that is independent of  $\mathbf{h}$ . The second term is minimized at 0 when  $\mathbf{h}(\mathbf{x}) = \mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\boldsymbol{\gamma}(\mathbf{z})]$ , which in turn minimizes (5).  $\square$

Therefore, minimizing (5) effectively minimizes the second term of (15) which is the expected l-2 squared distance between the prediction  $\mathbf{h}(\mathbf{x})$  and the DDC of exact posterior under  $\boldsymbol{\gamma}(\mathbf{z})$ , which depends on the flexibility of  $\mathbf{h}$ .

### A.2 MSE and the expected KL divergence

We first review a few known results for minimal exponential family from [46].

**Definition 1.** (Minimal exponential family [46, Section 3.2]) *A minimal exponential family distribution has the form*

$$q(\mathbf{z}) = \exp(\boldsymbol{\theta} \cdot \boldsymbol{\gamma}(\mathbf{z}) - \Phi(\boldsymbol{\theta})) \quad (16)$$

*in which there does not exist a nonzero real vector  $\mathbf{a}$  such that the linear combination  $\mathbf{a} \cdot \boldsymbol{\gamma}(\mathbf{z})$  is equal to a constant.*

If  $\boldsymbol{\gamma}$  is chosen to be a nonlinearity on random linear projections of  $\mathbf{z}$ , e.g.  $\gamma_i = \tanh(\mathbf{v}_i \cdot \mathbf{z} + b)$  with elements of  $\mathbf{v}_i$  and  $b$  being draws from a random distribution, then the  $\boldsymbol{\gamma}$  is linearly independent with probability one.

**Lemma 1.** (Log normalizer derivatives [46, Proposition 3.1]) *Let  $\mathbf{r}_Z(\boldsymbol{\theta}) = \mathbb{E}[\boldsymbol{\gamma}(\mathbf{z})]$  be the mean parameter of a minimal exponential family distribution in (16), the following holds:*

$$\frac{\partial \Phi(\boldsymbol{\theta})}{\partial \theta_i} = r_{Z,i}(\boldsymbol{\theta}) = \mathbb{E}[\gamma_i(\mathbf{z})] \quad (17)$$

$$\frac{\partial^2 \Phi(\boldsymbol{\theta})}{\partial \theta_i \partial \theta_j} = \frac{\partial r_{Z,i}(\boldsymbol{\theta})}{\partial \theta_j} = \mathbb{E}[\gamma_i(\mathbf{z})\gamma_j(\mathbf{z})] - \mathbb{E}[\gamma_i(\mathbf{z})]\mathbb{E}[\gamma_j(\mathbf{z})] =: \mathbb{C}[\boldsymbol{\gamma}(\mathbf{z})]_{ij} \quad (18)$$

Note that  $\nabla_{\boldsymbol{\theta}} \Phi(\boldsymbol{\theta})$  maps from  $\boldsymbol{\theta}$  to  $\mathbf{r}$  if and only if the exponential family distribution is minimal [46, Proposition 3.2]. In addition, under the same condition, there exists a mapping  $\boldsymbol{\theta}(\mathbf{r})$  such that  $\mathbb{E}[\boldsymbol{\gamma}(\mathbf{z})] = \mathbf{r}$ . Thus, the exponential family defined by the sufficient statistics  $\boldsymbol{\gamma}(\mathbf{z})$  can be specified

by either  $\theta$  or  $r$ . Importantly,  $r$  is a valid or feasible mean parameter if there exists some  $q$  such that  $\mathbb{E}_q[\gamma(z)] = r$ . Thus,  $\gamma$  defines a family of distributions by the set of all feasible mean parameters.

Let an internal model take joint distribution  $p(z, x)$ . Given a posterior DDC  $r(x) = h_\phi(x)$ , let the implied (by maximum entropy) exponential family distribution be  $q_\phi(z|x) := \exp(\theta(r(x)) \cdot \gamma(z) - \Phi(\theta(r(x))))$ . Let the error between the predicted and true DDC for a given  $\phi$  be  $e_\phi(x) = h_\phi(x) - \mathbb{E}_{p(z|x)}[\gamma(z)] = r(x) - \mathbb{E}[\gamma(z)|x]$ .

**Theorem 1.** *Under the following assumptions:*

- $\gamma(z)$  forms a minimal exponential family;
- $r(x)$  is a valid expectation under  $q_\phi(z|x)$  for any  $x$  ( $r(x)$  is in the set of feasible means);

If  $e_\phi(x) = 0$  for some  $\phi^*$  for all  $x$ , then  $\nabla_\phi \text{KL}[p(z|x)||q_\phi(z|x)]|_{\phi^*} = 0$ . Further, using  $h_\phi(x) = \mathbf{W}\sigma(x)$  ( $\phi = \mathbf{W}$ ) as the recognition model, and let  $\mathbf{W}^*$  be the minimizer of the EMSE problem in (5). If there exists  $\epsilon_{\mathbf{W}^*} > 0$  such that  $\mathbb{E}_{p(z)}[\|e_{\mathbf{W}^*}(x)\|_2^2] \leq \epsilon_\phi^2$ , and there exists an order 3 tensor  $\mathbf{A}$  and  $\epsilon_c > 0$  such that  $\mathbb{E}_{p(x)}[\|e_c(x)\|_2^2] \leq \epsilon_c^2$  where  $e_c(x) = \nabla_{\mathbf{W}}\theta(r(x))|_{\mathbf{W}^*} - \mathbf{A}\sigma(x)$  then

$$\|\mathbb{E}_{p(x)}[\nabla_{\mathbf{W}} \text{KL}[p(z|x)||q_{\mathbf{W}}(z|x)]|_{\mathbf{W}=\mathbf{W}^*}]\|^2 \leq \epsilon_c \epsilon_\phi. \quad (19)$$

*Proof.* The proof uses the same technique to show that the Expectation Propagation algorithm with exponential family factors minimizes a similar KL. For brevity, let  $\text{KL}[p||q] := \text{KL}[p(z|x)||q_\phi(z|x)]$  (note that  $e_c$  is a matrix)

$$\begin{aligned} \text{KL}[p||q] &= \int p(z|x) [\log p(z|x) - \log q_\phi(z|x)] dz \\ &= - \int p(z|x) [\log q_\phi(z|x)] dz \\ \nabla_\phi \text{KL}[p||q] &= - \int p(z|x) [\nabla_\phi \theta(r) \gamma(z) - \nabla_\phi \Phi(\theta(r))] dz \\ &= - \int p(z|x) \left[ \nabla_\phi \theta(r) (\gamma(z) - \frac{d\Phi(\theta(r))}{d\theta(r)}) \right] dz \\ &= \nabla_\phi \theta(r) [\mathbb{E}[\gamma(z)|x] - r] \\ &= \nabla_\phi \theta(r) [e_\phi(x)]. \end{aligned} \quad (20)$$

The second to last equality follows (17). Clearly,  $\nabla_\phi \text{KL}[p||q] = 0$  if  $e_{\mathbf{W}}(x) = 0, \forall x$ .

Now suppose  $r(x) = \mathbf{W}\sigma(x)$ . We decompose  $\nabla_{\mathbf{W}}\theta(r)$  as follows

$$\nabla_{\mathbf{W}}\theta(r(x)) = \nabla_{\mathbf{W}}\theta(r) - \mathbf{A}r(x) + \mathbf{A}r(x) \quad (21)$$

$$= e_c(x) + \mathbf{A}r(x). \quad (22)$$

Substituting in (20) and taking the expectation over  $p(x)$  gives

$$\begin{aligned} \mathbb{E}[\nabla_{\mathbf{W}} \text{KL}[p||q]] &= \mathbb{E}[e_c(x)e_{\mathbf{W}}(x) + (\mathbf{A}r(x)) \cdot e_{\mathbf{W}}(x)] \\ &= \mathbb{E}[e_c(x)e_{\mathbf{W}}(x)] + \mathbf{A}\mathbf{W}\mathbb{E}[\sigma(x)e_{\mathbf{W}}^T(x)] \\ \mathbb{E}[\|\nabla_{\mathbf{W}} \text{KL}[p||q]|_{\mathbf{W}^*}\|] &\stackrel{(1)}{\leq} \sqrt{\mathbb{E}[\|e_c(x)\|_2^2]} \sqrt{\mathbb{E}[\|e_{\mathbf{W}^*}(x)\|_2^2]} + \|\mathbf{A}\|_2^2 \|\mathbf{W}^*\|_2^2 \mathbb{E}[\|\sigma(x)e_{\mathbf{W}^*}^T(x)\|_2^2] \\ &\stackrel{(2)}{\leq} \epsilon_c \epsilon_{\mathbf{W}^*}, \end{aligned}$$

where (1) is due to the Cauchy-Schwarz inequality, and (2) is because the last term is the gradient of the EMSE w.r.t.  $\mathbf{W}$ , which is 0 when using  $\mathbf{W}^*$  that solving the EMSE problem:

$$\mathbb{E}_x[\sigma(x)e_{\mathbf{W}^*}^T] = \mathbb{E}_x[\sigma(x)(\mathbb{E}_{z|x}\gamma(z) - \mathbf{W}^*\sigma(x))^T] = \mathbb{E}_{z,x}[\sigma(x)(\gamma(z) - \mathbf{W}^*\sigma(x))^T] = 0.$$

□

The first assumption holds almost always. The second assumption is in general hard to reinforce, but after optimizing  $\phi$ , the DDC  $r(x)$  is likely to be inside the set of feasible means unless the true

posterior is close to a delta distribution on a single value of  $z$ , in which case the true posterior mean lies close to the boundary of the feasible set, and estimation error is likely to push  $r$  out of the feasible set.

The bound in (19) suggests that whenever  $\epsilon_W$  is small, the gradient of the KL is also small. For finite independent samples from  $p$ ,  $\epsilon_W$  shrinks at rate  $1/\sqrt{n}$ . The multiplier  $\epsilon_c$  suggests that the gradient of KL goes to zero faster if  $\sigma(x)$  better approximates the Jacobian  $\theta(r)$  w.r.t  $r$  ( $\epsilon_c$  is small). This Jacobian is, after inverting the total derivative  $\frac{dr}{d\theta}$  and using (18), is  $[\mathbb{C}_{q(x|x)}(\gamma(x))]^{-1}$ , which depends on the exponential family defined by  $\gamma(z)$ .

Thus, Theorem 1 suggests that an ideal  $\sigma(x)$  would be rich enough to linearly approximate not just the posterior mean but also the posterior covariance of  $\gamma(z)$  for all  $x$ . A simple  $\gamma(z)$  would help a given  $\sigma(x)$  satisfy these requirements, but a too simple  $\gamma(z)$  may not be rich enough to approximate a more complicated distribution, and the lowest KL could still be large even after optimizing the recognition parameters.

## B Formal solution to the filtering loss

We show the formal solution to minimizing (10) before discussing its biological implications.

**Proposition 2.** *Given a DDC of previous belief  $r_{t-1}$ ,  $W_{r_{t-1}}$  below is the minimizer of (10)*

$$\mathcal{L}^f(W) = \mathbb{E}_{q(z_{1:t}, x_t | x_{1:t-1})} [\|W\sigma(x_t) - \psi(z_{1:t})\|_2^2] \quad (10 \text{ revisited})$$

$$\begin{aligned} W_{r_{t-1}} &= C_{Z_{1:t}, X_t | x_{1:t-1}} C_{X_t, X_t | x_{1:t-1}}^{-1} \\ C_{Z_{1:t}, X_t | x_{1:t-1}} &= C_{Z_{1:t}, X_t | Z_{t-1}} r_{t-1} \quad C_{X_t, X_t | x_{1:t-1}} = C_{X_t, X_t | Z_{t-1}} r_{t-1} \\ C_{Z_{1:t}, X_t | Z_{t-1}} &= \arg \min_C \mathbb{E}_{p(z_{t-1}, z_t, x_t)} \|C\psi_{t-1} - \psi(z_{1:t})\sigma(x_t)\|_2^2 \\ C_{X_t, X_t | Z_{t-1}} &= \arg \min_C \mathbb{E}_{p(z_{t-1}, x_t)} \|C\psi_{t-1} - \sigma(x_t)\sigma(x_t)^\top\|_2^2. \end{aligned} \quad (23)$$

This is similar to the kernel Bayes rule [18]. The two minimization problems are essentially computing the readout weights used to approximate the conditional covariance matrices  $C$ . This solution for filtering involves solving these two problems before taking an inverse of a correlation matrix. If one interprets the two tensor  $C$ 's as weights, the matrix  $C$ 's are readout from  $r_{t-1}$ , then it is not clear how the inverse and  $W_{r_{t-1}}$  could be implemented by neural mechanisms.

## C Approximated solution for filtering

### C.1 The bilinear approximation and the tensor train decomposition

The bilinear approximation  $h_W(r_{t-1}, x_t)$  (12) and the corresponding solution to minimizing the EMSE (11) w.r.t.  $W$  is connected to the tensor train decomposition (TT) [35]. The EMSE is

$$\mathcal{L}^{bil}(W) = \mathbb{E}_{q(z_{1:t}, x_t, x_{1:t-1})} [\|W \cdot (r_{t-1} \otimes x_t) - \psi(z_{1:t})\|_2^2]. \quad (24)$$

Denote the minimizer of (24) at each  $t$  by  $W_t^*$ . Consider the situation that, at each  $t$ , we would like to predict  $\psi_t$  using a sequence of observations  $x_{1:t}$ . Let  $\sigma(\cdot) \in \mathbb{R}^{K_\sigma}$  be sufficiently rich so that there exists a linear operator  $W_t^{(p)}$  that maps from the product space of  $\sigma(x_1) \otimes \cdots \otimes \sigma(x_t)$  to  $r_t := \mathbb{E}_q(z_{1:t} | x_{1:t})[\psi(z_{1:t})]$ , then  $W_t^{(p)}$  is an order  $t+1$  tensor which is expensive to estimate. Low rank approaches may alleviate the difficulty, such as TT. In fact, the sequence of minimizers to (24)

$\{\mathbf{W}_{t'}^*\}_{t'=1}^t$  form a TT of an order  $t + 1$  tensor  $\mathbf{W}_t^{(f)}$  with the same shape as  $\mathbf{W}_t^{(p)}$ . For example:

$$\begin{aligned}
\mathbf{W}_1^* &= \arg \min_{\mathbf{W}_1} \mathbb{E}_p \sum_i \left( \sum_j W_{1,ji} \sigma_i(\mathbf{x}_1) - \psi_{1,j} \right)^2 \Rightarrow r_{1,j} = \sum_{ji} \mathbf{W}_{1,ji}^* \sigma_i(\mathbf{x}_1) \\
\mathbf{W}_2^* &= \arg \min_{\mathbf{W}_2} \mathbb{E}_p \sum_l \left( \sum_{jk} W_{2,lkj} r_{1,j} \sigma_k(\mathbf{x}_2) - \psi_{2,l} \right)^2 \\
&= \arg \min_{\mathbf{W}_2} \mathbb{E}_p \sum_l \left( \sum_{jk} W_{2,lkj} \left[ \sum_{ij} W_{1,ji}^* \sigma_i(\mathbf{x}_1) \right] \sigma_k(\mathbf{x}_2) - \psi_{2,l} \right)^2 \\
&= \arg \min_{\mathbf{W}_2} \mathbb{E}_p \sum_l \left( \sum_{ik} \underbrace{\left[ \sum_j W_{2,lkj} W_{1,ji}^* \right]}_{W_{2,lki}^{(f)}} \sigma_i(\mathbf{x}_1) \sigma_k(\mathbf{x}_2) - \psi_{2,l} \right)^2
\end{aligned}$$

the summation in the square brackets is the TT of  $\mathbf{W}_2^{(f)}$ . Thus, the proposed optimization for (24) finds a tensor of the same shape as  $\mathbf{W}_t^{(p)}$  in the TT space sequentially, predicting a new  $\psi_t$  by joining a new core tensor with  $\mathbf{W}_{t-1}^{(f)}$ , and only minimize the EMSE in the space of the new core tensor to get  $\mathbf{W}_t^*$ .

We argue that the computed  $\mathbf{r}_t$  after sequentially optimizing  $\mathbf{W}$  given a large set of training examples containing bootstrapped  $\mathbf{r}_{t-1}$  does not diverge and produce a good approximation to the true posterior moments on  $\psi_t$ . For any  $t$ , the set of inputs in the regression contains  $\mathbf{x}_t$ , so if the regression is performed in closed-form rather than using the delta rule, the output  $\mathbf{r}_t$  is at least close to the true  $\mathbb{E}_{p(\mathbf{z}_{1:t}|\mathbf{x}_{1:t})}[\psi_t]$  as  $\mathbb{E}_{q(\mathbf{z}_{1:t}|\mathbf{x}_t)}[\psi_t]$ , the output of another regression similar to (5), which can be made closer to  $\mathbb{E}_{p(\mathbf{z}_{1:t}|\mathbf{x}_t)}[\psi_t]$  using more flexible  $\mathbf{h}$  and more training examples. At time  $t + 1$ , the statistical dependency between  $\mathbf{r}_t$  (depending on  $\mathbf{x}_t$ ) and  $\psi_{t+1}$  improves the prediction quality if  $\mathbf{h}$  is flexible enough to pick up this dependency. At time  $t + \tau$ , as  $\tau > 0$  increases, the prediction should continue to improve until  $\mathbf{x}_t$  become uninformative of  $\psi_{t+\tau}$ , which depends on the range of temporal dependencies (“time constant”) of the internal model and the encoding functions  $\psi_t$ .

## D Experimental details

In all simulations in the main text, we assume the brain can draw samples from the internal model, and the recognition weights  $\mathbf{W}$  and readout weights  $\alpha$  have been trained on these samples for a long time and have converged. In our experiments, this condition was achieved by closed-form regression in solving least square regressions, using 10,000-20,000 sequences from the internal model and a Tikhonov regularization on  $\mathbf{W}$  with strength 0.001, and trained the recognition parameters for around 100 time steps in order for the SSM to enter in the stationary regime. The learned parameters are then fixed for online inference. The base tuning functions  $\gamma(\cdot)$  in (9) and input feature map  $\sigma(\cdot)$  in (12) and (13) have tanh nonlinearity after fixed random linear projections; the weights and biases in the projection are randomly drawn from a Gaussian with variance such that these functions are relatively smooth for the inputs they receive. Code is available at [https://github.com/kevin-w-li/ddc\\_ssm](https://github.com/kevin-w-li/ddc_ssm)

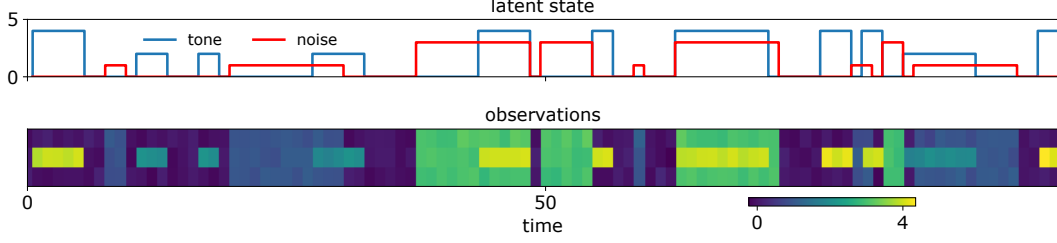


Figure 4: Example training data used for the auditory illusion experiment.

## D.1 Auditory continuity illusions

### D.1.1 Model setup

The internal model has a 2-D binary latent dynamics for the tone ( $z_{t,0}$ ) and noise ( $z_{t,1}$ ), and a 3-D noisy observation  $x_{t,i}, i \in \{0, 1, 2\}$  for three frequency bands. Mathematically, it is defined as

$$\begin{aligned}
 c_{t,i} &\sim \text{Bernoulli}(0.1) & i &\in \{0, 1\} \\
 l_{t,0} &\sim \text{Uniform}(\{2, 4\}) & l_{t,1} &\sim \text{Uniform}(\{1, 3\}) \\
 z_{t,i} &= \begin{cases} z_{t-1,i} & \text{if } z_{t-1,i} \neq 0.0 \text{ and } c_{t,i} = 0 \\ c_{t,i} l_{t,i} & \text{if } z_{t-1,i} = 0.0 \\ 0 & \text{if } z_{t-1,i} \neq 0.0 \text{ and } c_{t,i} = 1 \end{cases} & i &\in \{0, 2\} \\
 x_{t,1} &= \max\{z_{t,0}, z_{t,1}\} + \zeta_{t,i}, & \zeta_{t,1} &\sim \mathcal{N}(0, 0.1^2) \\
 x_{t,i} &= z_{t,1} + \zeta_{t,i}, & \zeta_{t,i} &\sim \mathcal{N}(0, 0.1^2) \quad i \in \{0, 2\}
 \end{aligned}$$

In words, the tone has energy levels  $\{0, 2, 4\}$  and the noise has energy levels  $\{0, 1, 3\}$ . At each time step, the tone and the noise can turn on or fall off with probability 0.1. For each of the two, if it turns on, it takes one of the two non-zero levels with equal chance; but it can only fall down to 0. The middle frequency channel reflects the greater level of the tone and the noise. The other two frequency channels only contain the noise. All three bands are contaminated by a small amount of i.i.d Gaussian noise. Example of the simulated data are shown in Figure 4.

In the DDC filter, we set the  $K_\psi = 200$ ,  $K_\gamma = 20$  and  $K_\sigma = 10$  and used the  $\mathbf{h}^{bil}$  in (12).

### D.1.2 Additional methods and results

We showed in Figure 1 the marginal p.m.f of the inferred tone level given observations up to time the stimulus time  $p(z_{t-\tau} | \mathbf{x}_{1:t})$  decoded from  $\mathbf{r}_t$ . This is done by first approximating posterior expectation over the static tuning function  $\gamma(z_{t-\tau})$  (other choices of basis are possible) using (14), obtaining  $\mathbf{m}_{t-\tau} := \mathbb{E}_{q(\mathbf{z}_{t-\tau} | \mathbf{x}_{1:t})} [\mathbf{Z}_{t-\tau} | \mathbf{x}_{1:t}]$ , a DDC on  $\mathbf{Z}_{t-\tau} | \mathbf{x}_{1:t}$ . Using maximum entropy decoding, we can find the corresponding p.m.f. Let the discrete p.m.f be  $p(z_{t-\tau} | \mathbf{x}_{1:t}) = \prod_i^{|\mathcal{Z}|} p_i^{\delta(z_{t-\tau} = z_i)}$ , where  $|\mathcal{Z}|$  is the cardinality of the support on  $\mathbf{z}$  (9 in this case), and  $\pi$  is the discrete probabilities that can be decoded from  $\mathbf{r}$  and  $\gamma$  by solving the following optimization problem:

$$\min_{\mathbf{p}} \sum_i^{|\mathcal{Z}|} p_i \log(p_i) \quad \text{s.t.} \quad \sum_i^{|\mathcal{Z}|} p_i \gamma_j(z_i) = m_j, \sum_i^{|\mathcal{Z}|} p_i = 1, p_i \in [0, 1], \quad (26)$$

which is relatively simple for a 9-outcome (3 tone  $\times$  3 noise levels), discrete distribution.

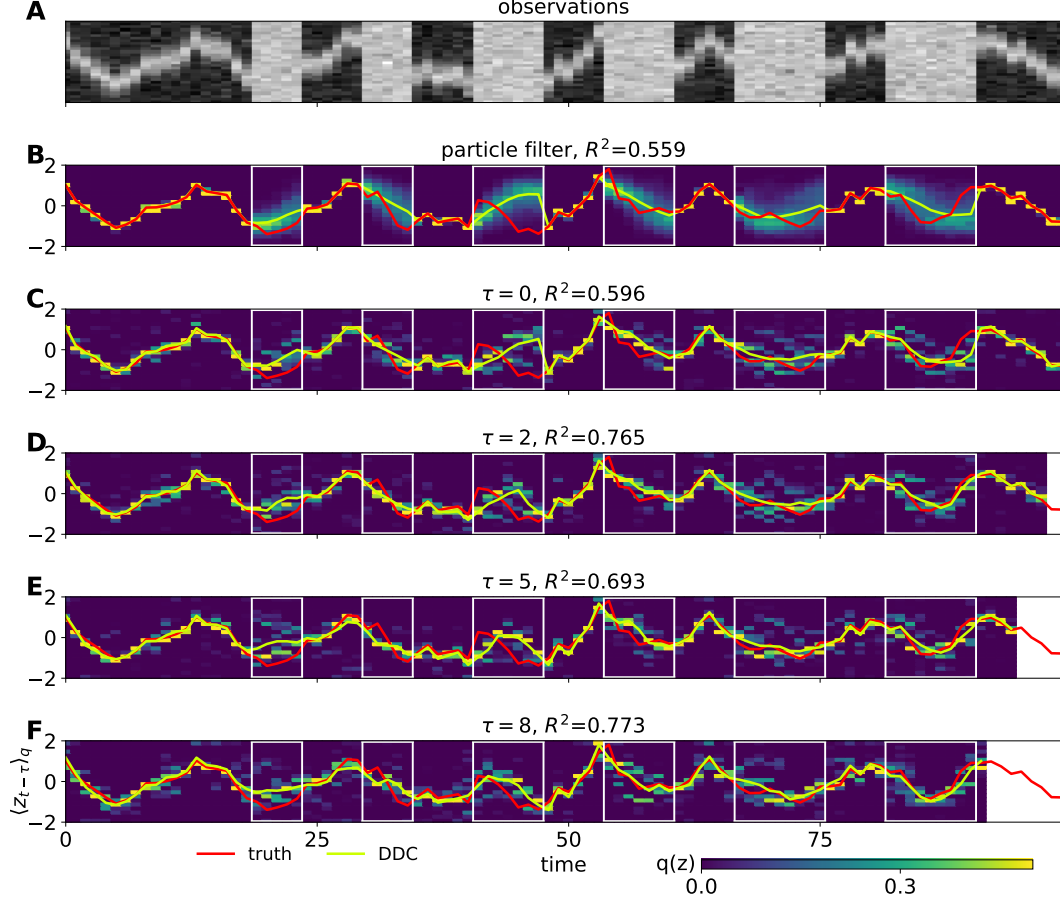


Figure 5: Maximum entropy decoding of the posterior marginals in the tracking experiment, compared with Figure 3 which is obtained by approximating expectation of bin functions.

## D.2 Flash-lag effect

The internal model that reproduced the smoothing effect is

$$p(\mathbf{z}_t | \mathbf{z}_{t-1}) = \mathcal{N}([A\mathbf{z}_{t-1}]_+, [0.01^2, 0.002^2, 1e^{-15}]) \quad A = \begin{bmatrix} 0.0 & 1.0 & 0.0 \\ 0.0 & 0.0 & 1.0 \\ 0.0 & 0.0 & 0.8 \end{bmatrix} \quad (27)$$

$$p(x_{t,i} | \mathbf{z}_t) = \text{Poisson} \left( 3 \exp \left[ -\frac{(\text{loc}(i) - z_{t,0})^2}{2 \times 1.5^2} \right] \right) \quad (28)$$

where  $[]_+$  is a elastic bounding box at  $\pm 1$ .  $\text{loc}$  is a linear transformation from pixel numbers to real values.

In the DDC filter, we set the dimensionalities  $K_\psi = 500$ ,  $K_\gamma = 100$  and  $K_\sigma = 150$  and used the  $\mathbf{h}^{\text{lin}}(\cdot)$  in (13).

## D.3 Noisy and occluded tracking

The internal model has 3-D latent (2 continuous, 1 discrete) and 30-D observation.

$$p(\mathbf{z}_t|\mathbf{z}_{t-1}) = \mathcal{N}(\mathbf{f}(\mathbf{z}_{t-1}), [0.1^2, 0.1^2]) \quad (29)$$

$$\mathbf{f}(\mathbf{z}_t) = s_t \mathbf{A} \mathbf{z}_{t-1} \quad (30)$$

$$s_t = \frac{1}{\|\mathbf{z}_{t-1}\|_2 \exp(-4(\|\mathbf{z}_{t-1}\|_2 - 0.3) + 1)} \quad (31)$$

$$p(m_t|m_{t-1}) = (\text{Bernoulli}(0.1) + m_{t-1}) \mod 2 \quad (32)$$

$$p(z_{t,i}|\mathbf{z}_t, m_t) = \mathcal{N}\left(\max\left\{\exp\left[-\frac{(\text{loc}(i) - z_{t,0})^2}{2 \times 3^2}\right], m_t\right\}, I_{30} 0.1^2\right) \quad (33)$$

$$(34)$$

where  $\text{loc}$  is a linear transformation from pixel number to real values, and  $\mathbf{A}$  is a rotation matrix by  $\pi/8$ . Due to the sigmoidal scaling,  $\mathbf{z}_t$  stays around the unit circle most of the time, but can occasionally cross through the origin due to noise.

In the DDC filter, we set the dimensionalities  $K_\psi = 500$ ,  $K_\gamma = 100$  and  $K_\sigma = 200$  and used the  $\mathbf{h}_{\mathbf{W}}^{\text{lin}}(\cdot)$  in (13).

The histogram decoding from  $\mathbf{r}_t$  is expected to be noisier due to the non-smoothness of the bin functions Figure 3, but still shows meaningful temporal integration of the observations. Results of the maximum entropy decoding of the posterior marginals are shown in Figure 5, which is less smooth due to  $\mathbf{r}_t$  being not exactly in the set of feasible sufficient statistics.

## E Robustness against neuronal noise

In the main text, we have discussed DDC when the representation  $\mathbf{r}_{\mathbf{Z}|\mathbf{x}}$  is deterministic, but real neurons are noisy. In this case, the spike count in some time window are taken as a noisy DDC representation. A noisy DDC may not identify a member in the class of exponential family distributions specified by  $\gamma$ , as it may not correspond to any valid mean parameter. However, if the noise has zero mean (including Poisson noise) and no or weak correlation, then it does not fatally harm inference or learning to infer, as long as the training for  $\mathbf{W}$  and  $\alpha$  is performed also on noisy DDC and on function evaluations on noisy samples. For inference, this type of noise on the input tends to average out in summations or inner products, the main operations in DDC computations as in (2) and (6). For learning to infer, noise on the input acts as a regularizer for  $\mathbf{W}$ , and noise on the output does not change the solution to regressions.

To verify our intuitions, we re-ran the experiments with the following changes to the DDC filter:

- The nonlinearity in  $\gamma(\mathbf{z}_t)$  and  $\sigma(\mathbf{x}_t)$  changes from  $\tanh(\cdot)$  to  $\text{sigmoid}(\cdot)$
- Independent Poisson noise is added to each feature evaluation of  $\gamma(\mathbf{z}_t)$ ,  $\sigma(\mathbf{x}_t)$  and  $\mathbf{k}(\mathbf{r}_{t-1}, \mathbf{x}_t)$  and output of  $\mathbf{h}_\phi$ .

and the results are shown in Figure 6 for the smoothing in flash-lag effect and Figure 7 for occluded tracking. The results are mostly the same as using noiseless DDC, although higher variability in prediction resulting from a noisy representation is clearly visible.

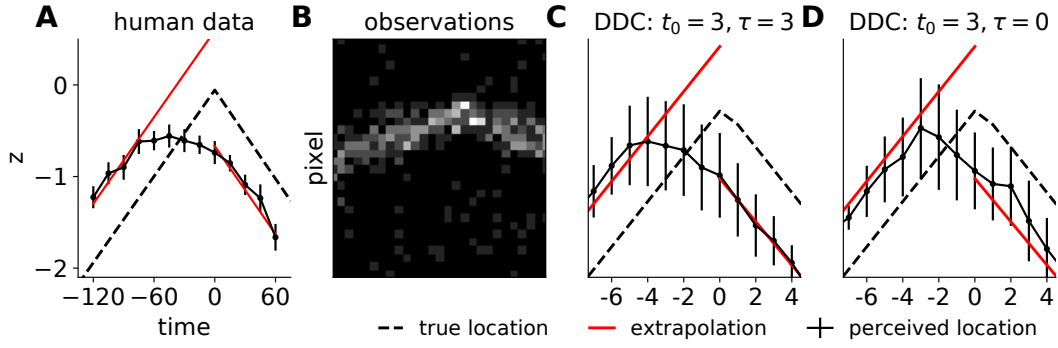


Figure 6: Same as Figure 2 but with noisy DDC. The error bars of DDC models are stds from 100 runs.

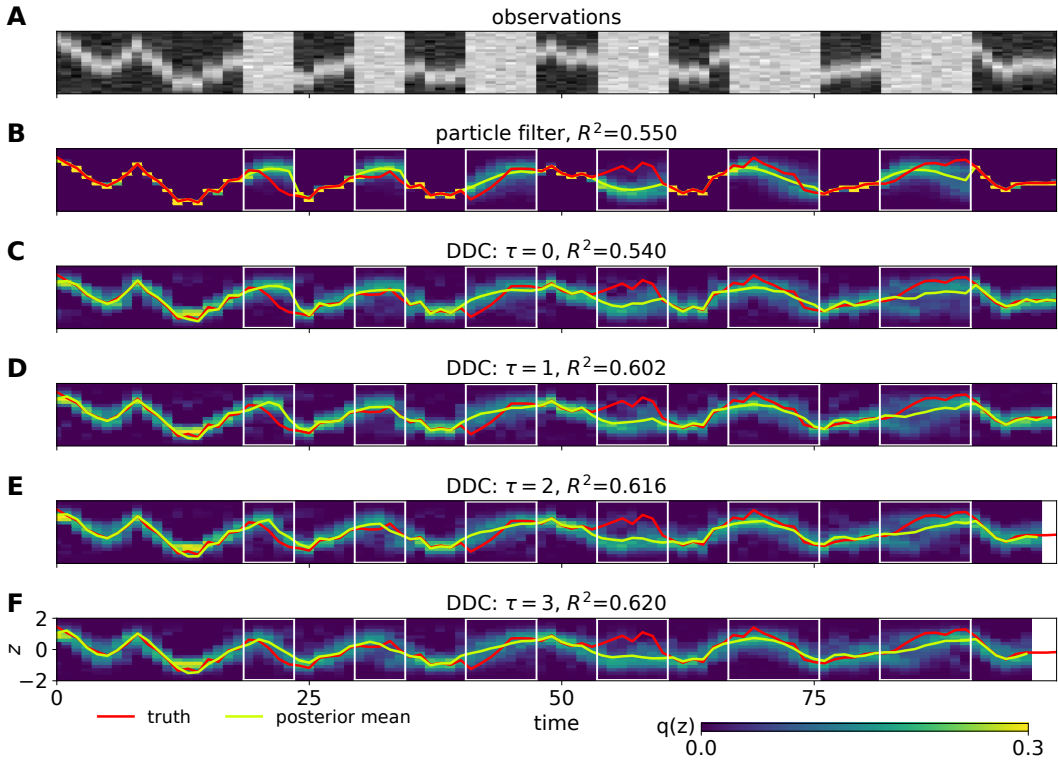


Figure 7: Same as Figure 3 but with noisy DDC.