# Appendix

This appendix describes the Bayesian Laplace propagation algorithm we derived for the two variants of the hierarchical model we presented in the main text. First, we describe Model I for capturing nonstationarity in firing rates, and then we move to Model II for capturing nonstationarity in neural dynamics.



Figure 1. Schematic of hierarchical non-stationary Poisson observation Latent Dynamical System (N-PLDS) models. Model I for capturing non-stationarity in mean firing rates. The parameter h slowly varies across trials and leads to fluctuations in mean firing rates. Model II for capturing non-stationarity in population dynamics. The dynamics matrix A changes across trials, as controlled by the hyperparameters Φ.

## Model I : nonstationarity in firing rates

# Basic setup

Likelihood:  $\mathbf{y}_t \in \mathbb{R}^p$ ,  $\mathbf{x}_t \in \mathbb{R}^k$ ,  $C \in \mathbb{R}^{p \times k}$ 

$$
p(\mathbf{y}_t|\mathbf{x}_t, C, \mathbf{h}^{(i)}) = \text{Poiss}(\mathbf{y}_t|\exp(C(\mathbf{x}_t + \mathbf{h}^{(i)}) + \mathbf{d})),
$$

where  $\mathbf{h}^{(i)} \in \mathbb{R}^k$  is a vector of latent variables that capture nonstationarity in firing rates across recordings  $i = \{1, \dots, r\}$ . Latent dynamics:  $A \in \mathbb{R}^{k \times k}$  and  $B \in \mathbb{R}^{k \times d}$ 

$$
p(\mathbf{x}_t|\mathbf{x}_{t-1}, A) = \mathcal{N}(\mathbf{x}_t|A\mathbf{x}_{t-1} + B\mathbf{u}_t, I).
$$

Parameters in this model:  $\Theta = \{A, B, C, \mathbf{h}^{(1:r)}\}$ . For simplicity, we set **d** to its ML estimate. Vectorized notations:  $\mathbf{a} = \text{vec}(A^{\top}) \in \mathbb{R}^{k^2}, \, \mathbf{b} = \text{vec}(B^{\top}) \in \mathbb{R}^{kd}, \, \text{and} \, \mathbf{c} = \text{vec}(C^{\top}) \in \mathbb{R}^{pk}.$ 

Priors:

$$
p(\mathbf{a}|\alpha) = \mathcal{N}(\mathbf{a}|\mathbf{0}, \alpha^{-1}\mathbf{I}), \quad \mathbf{p}(\mathbf{b}|\beta) = \mathcal{N}(\mathbf{b}|\mathbf{0}, \beta^{-1}\mathbf{I}), \tag{1}
$$

For  $\mathbf{h}^{(i)}$ , we assume slowly varying dynamics across recordings

$$
\mathbf{h}^{(i)} \sim \mathcal{GP}(\mathbf{m}_{\mathbf{h}}, K(i,j)) \tag{2}
$$

where we denote the (vector) mean and (matrix) covariance functions by  $K(i, j)$ , respectively, where the  $(i, j)$ th block of the covariance matrix is given by

$$
K(i,j) = (\sigma^2 + \epsilon \delta_{i,j}) \exp\left(-\frac{1}{2\tau^2}(i-j)^2\right) I_k
$$
\n(3)

The hyperparameters in total are  $\Phi = {\mathbf{m_h}, \alpha, \beta, \sigma^2, \tau^2}.$ 

# Variational lower bound

The marginal likelihood of the observations is lower bounded by

$$
\log p(\mathbf{y}_{1:T}^{(1:r)}) \geq \int d\theta \, d\mathbf{x}_{1:T}^{(1:r)} \, q(\theta, \mathbf{x}_{1:T}^{(1:r)}) \, \log \frac{p(\theta, \mathbf{x}_{1:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)})}{q(\theta, \mathbf{x}_{1:T}^{(1:r)})},\tag{4}
$$

where the approximate posterior factor is

$$
q(\theta, \mathbf{x}_{1:T}^{(1:r)}) = q_{\theta}(\theta) \prod_{i=1}^{r} q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)}),
$$
\n
$$
(5)
$$

where 
$$
q_{\theta}(\theta) = q_{\mathbf{a},\mathbf{b}}(\mathbf{a},\mathbf{b}) q_{\mathbf{c},\mathbf{h}}(\mathbf{c},\mathbf{h}^{(1:r)})
$$
. (6)

For simplicity, we further assume

$$
q_{\mathbf{c},\mathbf{h}}(\mathbf{c},\mathbf{h}^{(1:r)}) = q_{\mathbf{h}|\mathbf{c}}(\mathbf{h}^{(1:r)}|\mathbf{c})q(\mathbf{c}), \qquad (7)
$$

$$
= q_{\mathbf{h}|\mathbf{c}}(\mathbf{h}^{(1:r)}|\mathbf{c})\delta(\mathbf{c}-\hat{\mathbf{c}}), \tag{8}
$$

$$
= q_{\mathbf{h}|\hat{\mathbf{c}}}(\mathbf{h}^{(1:r)}|\hat{\mathbf{c}}),\tag{9}
$$

where  $\hat{\mathbf{c}}$  is maximum likelihood estimate of  $\mathbf{c}$ .

# Bayesian Laplace propagation

### Posterior over parameters

We compute  $q_{\theta}(\theta)$  by integrating out latent variables from the total log joint distribution:

$$
\log q_{\theta}(\theta) = \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(1:r)})}\left[\log p(\mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)}, \theta)\right] + const,
$$
\n
$$
= \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)})}\left[\log p(\mathbf{y}_{1:T}^{(1:r)}|\mathbf{x}_{0:T}^{(1:r)}, \theta) + \log p(\mathbf{x}_{0:T}^{(1:r)}|\theta) + \log p(\theta)\right] + const,
$$
\n
$$
= \sum_{i=1}^{r} [\mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)})}(\sum_{t=1}^{T} (\log p(\mathbf{y}_{t}^{(i)}|\mathbf{x}_{t}^{(i)}, \mathbf{c}, \mathbf{h}^{(i)}) + \log p(\mathbf{x}_{t}^{(i)}|\mathbf{x}_{t-1}^{(i)}, \mathbf{u}_{t}, \mathbf{a}, \mathbf{b}))] + \log p(\mathbf{a}|\alpha) + \log p(\mathbf{b}|\beta) + \log p(\mathbf{h}^{(1:r)}|0, K) + const.
$$
\n(10)

Note that we assume the inputs **u** are the same across recordings. (so we don't put the recording index i on  $\mathbf{u}_t$ ).

#### 1. approximate posterior over a, b

We can compute  $q_{a,b}(a,b)$  by extracting all the terms in  $\log p(\mathbf{x}_{0:T}, \mathbf{y}_{1:T}, \theta)$  that depend on  $a, b$  and then taking the expectation of the terms w.r.t.  $q_{\mathbf{x}}(\mathbf{x}_{0:T})$ :

$$
\log q_{\mathbf{a},\mathbf{b}}(\mathbf{a},\mathbf{b}) = \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)})} \left[ \sum_{t=1}^{T} \log p(\mathbf{x}_{t}^{(i)} | \mathbf{x}_{t-1}^{(i)}, \mathbf{u}_{t}, \mathbf{a}, \mathbf{b}) \right] + \log p(\mathbf{a}|\alpha) + \log p(\mathbf{b}|\beta) + const,
$$
  
\n
$$
= -\frac{1}{2} \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)})} \left[ \sum_{t=1}^{T} (\mathbf{x}_{t}^{(i)} - A\mathbf{x}_{t-1}^{(i)} - B\mathbf{u}_{t})^{\top} (\mathbf{x}_{t}^{(i)} - A\mathbf{x}_{t-1}^{(i)} - B\mathbf{u}_{t}) \right] - \frac{\alpha}{2} \mathbf{a}^{\top} \mathbf{a} - \frac{\beta}{2} \mathbf{b}^{\top} \mathbf{b} + const,
$$
(11)  
\n
$$
= -\frac{1}{2} \sum_{i=1}^{r} \left[ \mathbf{a}^{\top} (\frac{\alpha}{r} I + W_{A^{(i)}}^{\mathbf{b}}) \mathbf{a} - 2\mathbf{a}^{\top} (\text{vec}(S_{A^{(i)}}) - G_{A^{(i)}}^{\mathbf{b}} \mathbf{b}) + \mathbf{b}^{\top} (\frac{\beta}{r} I + \ddot{U}^{\mathbf{b}d}) \mathbf{b} - 2\mathbf{b}^{\top} \text{vec}(\tilde{M}^{(i)}) \right] + const,
$$

where  $W_{A^{(i)}}^{bd} = I_r \otimes W_{A^{(i)}}, G_{A^{(i)}}^{bd} = I_r \otimes G_{A^{(i)}}, \ddot{U}^{bd} = I_r \otimes \ddot{U}$ , and the sufficient statistics are denoted by

$$
W_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \mathbf{x}_{t-1}^{(i)} | ^{\top} \rangle, \quad S_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \mathbf{x}_{t}^{(i)} | ^{\top} \rangle, \qquad G_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \rangle \mathbf{u}_{t}^{(i)} ,
$$

$$
\ddot{U} = \sum_{t=1}^{T} \mathbf{u}_{t} \mathbf{u}_{t}^{(i)} , \qquad \tilde{M}^{(i)} = \sum_{t=1}^{T} \mathbf{u}_{t} \langle \mathbf{x}_{t}^{(i)} \rangle ^{\top} \qquad (12)
$$

Using new notations  $\mathbf{W} = \sum_{i=1}^r W_{A^{(i)}}^{bd}$ ,  $\mathbf{s} = \sum_{i=1}^r \text{vec}(S_{A^{(i)}})$ ,  $\mathbf{G} = \sum_{i=1}^r G_{A^{(i)}}^{bd}$ ,  $\mathbf{m} = \sum_{i=1}^r \text{vec}(\tilde{M}^{(i)})$ , we rewrite eq. 11, whose derivative expressions are given by

$$
\log q_{\mathbf{a},\mathbf{b}}(\mathbf{a},\mathbf{b}) = -\frac{1}{2} \left[ \mathbf{a}^\top (\alpha I + \mathbf{W}) \mathbf{a} - 2 \mathbf{a}^\top (\mathbf{s} - \mathbf{G} \mathbf{b}) + \mathbf{b}^\top (\beta I + r \ddot{U}^{bd}) \mathbf{b} - 2 \mathbf{b}^\top \mathbf{m} \right],\tag{13}
$$

$$
H_{\mathbf{a}} = -\frac{\partial}{\partial \mathbf{a} \mathbf{a}^{\top}} \log q_{\mathbf{a}, \mathbf{b}}(\mathbf{a}, \mathbf{b}) = (\alpha I + \mathbf{W}), \tag{14}
$$

$$
H_{\mathbf{ab}} = -\frac{\partial}{\partial \mathbf{ab}^{\top}} \log q_{\mathbf{a},\mathbf{b}}(\mathbf{a},\mathbf{b}) = \mathbf{G},\tag{15}
$$

$$
H_{\mathbf{b}} = -\frac{\partial}{\partial \mathbf{b} \mathbf{b}^{\top}} \log q_{\mathbf{a}, \mathbf{b}}(\mathbf{a}, \mathbf{b}) = (\beta I + r\ddot{U}^{bd}), \tag{16}
$$

$$
\frac{\partial}{\partial \mathbf{a}} \log q_{\mathbf{a},\mathbf{b}}(\mathbf{a},\mathbf{b}) = -(\alpha I + \mathbf{W})\mathbf{a} + (\mathbf{s} - \mathbf{G}\mathbf{b}) = -H_{\mathbf{a}}\mathbf{a} + (\mathbf{s} - \mathbf{G}\mathbf{b}),\tag{17}
$$

$$
\frac{\partial}{\partial \mathbf{b}} \log q_{\mathbf{a},\mathbf{b}}(\mathbf{a},\mathbf{b}) = -\mathbf{G}^{\top} \mathbf{a} - (\beta I + r\ddot{U}^{bd})\mathbf{b} + \mathbf{m} = -\mathbf{G}^{\top} \mathbf{a} - H_{\mathbf{b}}\mathbf{b} + \mathbf{m}.
$$
 (18)

Using Schur complement, we obtain the covariance of  $q(\mathbf{a}, \mathbf{b})$ 

$$
\Sigma_{\mathbf{a}} = (H_{\mathbf{a}} - H_{\mathbf{a}\mathbf{b}} H_{\mathbf{b}}^{-1} H_{\mathbf{a}\mathbf{b}}^\top)^{-1},\tag{19}
$$

$$
\Sigma_{\mathbf{b}} = (H_{\mathbf{b}} - H_{\mathbf{a}\mathbf{b}}^\top H_{\mathbf{a}}^{-1} H_{\mathbf{a}\mathbf{b}})^{-1}, \tag{20}
$$

$$
\Sigma_{\mathbf{ab}} = -\Sigma_{\mathbf{a}} H_{\mathbf{ab}} H_{\mathbf{b}}^{-1},\tag{21}
$$

and the mean of  $q(\mathbf{a}, \mathbf{b})$ ,

$$
\mu_{\mathbf{b}} = \Sigma_{\mathbf{b}} (\mathbf{m} - H_{\mathbf{a}\mathbf{b}}^\top H_{\mathbf{a}}^{-1} \mathbf{s}), \tag{22}
$$

$$
\mu_{\mathbf{a}} = \Sigma_{\mathbf{a}} (\mathbf{s} - H_{\mathbf{a}\mathbf{b}} H_{\mathbf{b}}^{-1} \mathbf{m}). \tag{23}
$$

## 2. Computing  $q_{\mathbf{h}|\hat{\mathbf{c}}}(\mathbf{h}^{(1:r)}|\hat{\mathbf{c}})$

Assuming we have the maximum likelihood estimate of c, we write down all the terms in  $\log p(\mathbf{x}_{0:T}, \mathbf{y}_{1:T}, \theta)$  that depend on  $\mathbf{h}^{(1:r)}$ :

$$
\log q_{\mathbf{h}|\hat{\mathbf{c}}}(\mathbf{h}^{(1:r)}|\hat{\mathbf{c}}) = \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \left[ \sum_{t=1}^{T} \log p(\mathbf{y}_{t}^{(i)}|\mathbf{x}_{t}^{(i)}, \hat{\mathbf{c}}, \mathbf{h}^{(i)}) \right] + \log p(\mathbf{h}^{(1:r)}|0, K),
$$
  
\n
$$
= \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \left[ \sum_{t=1}^{T} (\mathbf{y}_{t}^{(i)\top}(\hat{C}(\mathbf{x}_{t}^{(i)} + \mathbf{h}^{(i)}) + \mathbf{d}) - \mathbf{1}^{\top} \exp(\hat{C}(\mathbf{x}_{t}^{(i)} + \mathbf{h}^{(i)}) + \mathbf{d})) \right] - \frac{1}{2} \mathbf{h}^{(1:r)\top} K^{-1} \mathbf{h}^{(1:r)} + const,
$$
  
\n
$$
= \sum_{i=1}^{r} \left[ \hat{C} S_{C^{(i)}} + \sum_{t=1}^{T} \mathbf{y}_{t}^{(i)\top} (\hat{C} \mathbf{h}^{(i)} + \mathbf{d}) - \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \mathbf{1}^{\top} \exp(\hat{C}(\mathbf{x}_{t}^{(i)} + \mathbf{h}^{(i)}) + \mathbf{d})) \right] - \frac{1}{2} \mathbf{h}^{(1:r)\top} K^{-1} \mathbf{h}^{(1:r)} + const,
$$
\n(24)

where each row of  $\hat{C}$  is denoted by  $\hat{\mathbf{c}}_s$  and the sufficient statistic is denoted by

$$
S_{C^{(i)}} = \sum_{t=1}^{T} < \mathbf{x}_t^{(i)} > \mathbf{y}_t^{(i)\top} \tag{25}
$$

Assuming the approximate posterior over latent variables is multivariate Gaussian with marginals  $q(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t | \boldsymbol{\omega}_t, \Upsilon_t)$ , the expectation of the exponential term above is given by

$$
\mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})}\left\{\sum_{t=1}^{T} \exp(\hat{\mathbf{c}}_{s}^{\top} \mathbf{x}_{t}^{(i)})\right\} = \int d\mathbf{x}_{1:T}^{(i)} q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)}) \exp(\hat{\mathbf{c}}_{s}^{\top} \mathbf{x}_{1}^{(i)} + \dots + \hat{\mathbf{c}}_{s}^{\top} \mathbf{x}_{T}^{(i)}),
$$
\n
$$
= \sum_{t=1}^{T} \exp(\hat{\mathbf{c}}_{s}^{\top} \boldsymbol{\omega}_{t}^{(i)} + \frac{1}{2}\hat{\mathbf{c}}_{s}^{\top} \boldsymbol{\Upsilon}_{t}^{(i)}\hat{\mathbf{c}}_{s}). \tag{26}
$$

Therefore, the log joint distribution is given by

$$
\log q_{\mathbf{h}|\hat{\mathbf{c}}}(\mathbf{h}^{(1:r)}|\hat{\mathbf{c}}) = \sum_{i=1}^{r} \left[ \hat{C}S_{C^{(i)}} + \sum_{t=1}^{T} \left( \mathbf{y}_{t}^{(i)\top}(\hat{C}\mathbf{h}^{(i)} + \mathbf{d}) - \mathbf{1}^{\top} \exp(\hat{C}(\omega_{t}^{(i)} + \mathbf{h}^{(i)}) + \frac{1}{2} \text{diag}(\hat{C}\Upsilon_{t}^{(i)}\hat{C}^{\top}) + \mathbf{d}) \right) \right] - \frac{1}{2}\mathbf{h}^{(1:r)\top}K^{-1}\mathbf{h}^{(1:r)},
$$
\n
$$
= \mathbf{h}^{(1:r)\top}(\hat{C}^{bd\top}\mathbf{y}^{(1:r)}) - \mathbf{1}^{\top}(\exp(\hat{C}^{bd}\mathbf{h}^{(1:r)}) \circ \mathbf{g}^{(1:r)}) - \frac{1}{2}\mathbf{h}^{(1:r)\top}K^{-1}\mathbf{h}^{(1:r)},
$$
\n(27)

 $\text{where } \hat{C}^{bd} = I_r \otimes \hat{C}, \ \mathbf{y}^{(1:r)} = [\sum_{t=1}^T \mathbf{y}_t^{(1)}, \cdots, \sum_{t=1}^T \mathbf{y}_t^{(r)}]^\top, \ \mathbf{g}^{(1:r)} = [\mathbf{g}^{(1)}, \cdots, \mathbf{g}^{(r)}]^\top, \ \text{where } \mathbf{g}^{(i)} = \sum_{t=1}^T \exp(\hat{C}\boldsymbol{\omega}_t^{(i)} + \hat{C})$  $\frac{1}{2}$ diag $(\hat{C}\Upsilon_t^{(i)}\hat{C}^\top)+\mathbf{d}$ ).

We approximate the joint posterior  $q_{\mathbf{h}|\hat{\mathbf{c}}}(\mathbf{h}^{(1:r)}|\hat{\mathbf{c}})$  as a Gaussian distribution from the derivatives w.r.t.  $\mathbf{h}^{(1:r)}$ :

$$
q_{\mathbf{h}|\hat{\mathbf{c}}}(\mathbf{h}^{(1:r)}|\hat{\mathbf{c}}) = \mathcal{N}(\mathbf{h}^{(1:r)}|\boldsymbol{\mu}_{\mathbf{h}}, \boldsymbol{\Sigma}_{\mathbf{h}})
$$
(28)

$$
\Sigma_{\mathbf{h}}^{-1} = H_{\mathbf{h}} + K^{-1}, \quad \text{where } \mathbf{h}^{(1:r)} = \boldsymbol{\mu}_{\mathbf{h}}, \tag{29}
$$

$$
\boldsymbol{\mu}_{\mathbf{h}} = K \left[ \hat{C}^{bd\top} \mathbf{y}^{(1:r)} - \hat{C}^{bd\top} (\exp(\hat{C}^{bd} \mathbf{h}^{(1:r)}) \circ \mathbf{g}^{(1:r)}) \right], \text{ where } \mathbf{h}^{(1:r)} = \boldsymbol{\mu}_{\mathbf{h}} \tag{30}
$$

where  $H_{\mathbf{h}} = -\frac{\partial^2}{\partial^2 \mathbf{h}^{(1)}}$  $\frac{\partial^2}{\partial^2 \mathbf{h}^{(1:h)}} \sum_{i=1}^r [\int d \mathbf{x}_{0:1}^{(i)}$  $_{0:T}^{\left( i\right) }q(\mathbf{x}_{0:T}^{\left( i\right) }% )=\left( \begin{smallmatrix} 0&1\\0&1\end{smallmatrix}\right) ^{T}q(\mathbf{x}_{0:T}^{\left( i\right) }% )$  $\sum_{t=1}^{(i)} \log p(\mathbf{y}_{t}^{(i)}|\mathbf{x}_{t}^{(i)}, \hat{\mathbf{c}}, \hat{\mathbf{d}}, \mathbf{h}^{(i)})] = \hat{C}^{bd\top} \text{diag}\left[\exp(\hat{C}^{bd}\mathbf{h}^{(1:r)})\circ \mathbf{g}^{(1:r)}\right] \hat{C}^{bd}.$ 

### 3. Computing the ML estimate of  $\hat{c}$

We set C to the ML estimate  $\hat{C}$ , which is obtained by

$$
\hat{\mathbf{c}} = \arg \max_{\mathbf{c}} \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \left[ \sum_{t=1}^{T} \log p(\mathbf{y}_t^{(i)} | \mathbf{x}_t^{(i)}, \mathbf{c}, \mathbf{h}^{(i)}) \right],
$$
\n(31)

whose first derivatives w.r.t.  $C$  is given by :

$$
\sum_{i=1}^{r} \left[ S_{C^{(i)}} \top + \sum_{t=1}^{T} (\mathbf{y}_t \mathbf{h}^{(i)\top} - l^{(i)}(C) (\boldsymbol{\omega}_t^{(i)} + \mathbf{h}^{(i)})^\top - \text{diag}(l^{(i)}(C)) \Upsilon_t^{(i)} C) \right]
$$
(32)

where we fix  $\mathbf{h}^{(i)}$  to its posterior mean  $\boldsymbol{\mu}_{\mathbf{h}^{(i)}}$  and  $l^{(i)}(C) = \exp(C^{\top}(\boldsymbol{\omega}_t^{(i)} + \mathbf{h}^{(i)}) + \frac{1}{2} \text{diag}(C^{\top} \Upsilon_t^{(i)} C) + \mathbf{d}).$ 

### 4. ML estimate of d

$$
\hat{\mathbf{d}} = \arg \max_{\mathbf{d}} \sum_{i=1}^{r} \left[ \hat{C} S_{C^{(i)}} + \sum_{t=1}^{T} \left( \mathbf{y}_{t}^{(i)\top} (\hat{C} \mathbf{h}^{(i)} + \mathbf{d}) - \mathbf{1}^{\top} \exp(\hat{C}(\boldsymbol{\omega}_{t}^{(i)} + \mathbf{h}^{(i)}) + \frac{1}{2} \text{diag}(\hat{C} \Upsilon_{t}^{(i)} \hat{C}^{\top}) + \mathbf{d}) \right) \right], \tag{33}
$$

$$
= \log(\sum_{i=1}^{r} \sum_{t=1}^{T} \mathbf{y}_{t}^{(i)}) - \log(\sum_{i=1}^{r} \sum_{t=1}^{T} \exp(\hat{C}(\boldsymbol{\omega}_{t}^{(i)} + \mathbf{h}^{(i)}) + \frac{1}{2} \text{diag}(\hat{C} \Upsilon_{t}^{(i)} \hat{C}^{\top}))).
$$
\n(34)

,

## Posterior over latent variables

In VBE step, we compute  $q_{\mathbf{x}}(\mathbf{x}_{0:qT})$  by

$$
\sum_{i=1}^{r} \log q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)}) = \sum_{i=1}^{r} \mathbb{E}_{q_{\theta}(\theta)} \log p(\theta, \mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)}) + const,\n= \sum_{i=1}^{r} \left[ \mathbb{E}_{q_{\theta}(\theta)} \log p(\mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)} | \theta) - \log Z'_{(i)} \right],
$$
\n(35)

where the normalization constant is given by

$$
Z'_{(i)} = \int d\mathbf{x}_{0:T}^{(i)} \exp\left(\mathbb{E}_{q_{\theta}(\theta)} \log p(\mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)} | \theta)\right).
$$
 (36)

The complete-data log likelihood in the ith recording is written as

$$
\log p(\mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)} | \theta) = \sum_{t=1}^{T} \{ \log p(\mathbf{y}_t^{(i)} | \mathbf{x}_t^{(i)}, C, \mathbf{d}, \mathbf{h}^{(i)}) + \log p(\mathbf{x}_t^{(i)} | \mathbf{x}_{t-1}^{(i)}, A, B, \mathbf{u}_t) \},
$$
(37)

which tells us that the log posterior over latent variables is quadratic in each  $x_t$ . This enables us to use the sequential update of the posterior over latent variables. We will also use the following sequential forward/backward algorithm for each recording in parallel. In the following, the recording index  $i$  on  $x, y$  is removed for notational cleanness.

### Forward filtering

We denote the posterior over the latent variables at each time  $t$  by

$$
\alpha(\mathbf{x}_t) \propto \int d\mathbf{x}_{t-1} \alpha(\mathbf{x}_{t-1}) \exp \left[ \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_\theta(\theta)} \right], \tag{38}
$$

$$
\propto \exp(\langle \log p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_\theta(\theta)} \left\{ \int d\mathbf{x}_{t-1} \alpha(\mathbf{x}_{t-1}) \exp\left(\langle \log (p(\mathbf{x}_t|\mathbf{x}_{t-1}) \rangle_{q(\theta)}) \right) \right\}.
$$
 (39)

Assuming  $\alpha(\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{t-1}, \Sigma_{t-1})$ , the integral is analytically tractable since the second part in the integrand is also quadratic in  $\mathbf{x}_{t-1}$ :

$$
\exp[-\frac{1}{2}(\mathbf{x}_{t-1}^\top < A^\top A > \mathbf{x}_{t-1} - 2\mathbf{x}_{t-1}^\top(< A > \mathbf{x}_{t-1}^\top < A^\top B > \mathbf{u}_t) + {\mathbf{x}_t}^\top \mathbf{x}_t - 2{\mathbf{x}_t}^\top < B > \mathbf{u}_t + {\mathbf{u}_t}^\top < B^\top B > \mathbf{u}_t)].
$$

The integrand is summarised as

$$
\alpha(\mathbf{x}_{t-1}) \exp\left(\langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) >_{q(\theta)}\right) = Z\mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{t-1}^*, \boldsymbol{\Sigma}_{t-1}^*),\tag{40}
$$

$$
\Sigma_{t-1}^{*-1} = \Sigma_{t-1}^{-1} + \langle A^{\top} A \rangle, \tag{41}
$$

$$
\mu_{t-1}^* = \Sigma_{t-1}^*(\Sigma_{t-1}^{-1}\mu_{t-1} + ^\top \mathbf{x}\_t -  \mathbf{u}\_t\), \tag{42}
$$

and the remaining term  $Z$  is given by:

$$
Z = \exp[-\frac{1}{2}(\mathbf{x}_t^\top \mathbf{x}_t - 2\mathbf{x}_t^\top < B > \mathbf{u}_t + \mathbf{u}_t^\top < B^\top B > \mathbf{u}_t) + \frac{1}{2}\boldsymbol{\mu}_{t-1}^*^\top \boldsymbol{\Sigma}_{t-1}^{*-1} \boldsymbol{\mu}_{t-1}^*],
$$
\n(43)

where

$$
\frac{1}{2}\mu_{t-1}^* \nabla_{t-1}^{* - 1} \mu_{t-1}^* = \frac{1}{2} (\Sigma_{t-1}^{-1} \mu_{t-1} + \langle A \rangle^\top \mathbf{x}_t - \langle A^\top B \rangle \mathbf{u}_t)^\top \Sigma_{t-1}^* (\Sigma_{t-1}^{-1} \mu_{t-1} + \langle A \rangle^\top \mathbf{x}_t - \langle A^\top B \rangle \mathbf{u}_t),
$$
\n
$$
= \frac{1}{2} (\mathbf{x}_t^\top \langle A \rangle \Sigma_{t-1}^* \langle A \rangle^\top \mathbf{x}_t + 2 \mathbf{x}_t^\top \langle A \rangle (\Sigma_{t-1}^* \Sigma_{t-1}^{-1} \mu_{t-1} - \Sigma_{t-1}^* \langle A^\top B \rangle \mathbf{u}_t) +
$$
\n
$$
\mu_{t-1}^\top \Sigma_{t-1}^{-1} \Sigma_{t-1}^* \Sigma_{t-1}^{-1} \mu_{t-1} - 2 \mu_{t-1}^\top \Sigma_{t-1}^{-1} \Sigma_{t-1}^* \langle A^\top B \rangle \mathbf{u}_t + \mathbf{u}_t^\top \langle A^\top B \rangle^\top \Sigma_{t-1}^* \langle A^\top B \rangle \mathbf{u}_t).
$$

$$
Z \propto \mathcal{N}(\mathbf{x}_t | \tilde{\boldsymbol{\mu}}_t, \tilde{\Sigma}_t), \tag{44}
$$

$$
\tilde{\Sigma}_t^{-1} = I - \langle A \rangle \Sigma_{t-1}^* \langle A \rangle^\top,\tag{45}
$$

$$
\tilde{\boldsymbol{\mu}}_t = \tilde{\Sigma}_t \left( \langle B > \mathbf{u}_t + \langle A > \Sigma_{t-1}^* \Sigma_{t-1}^{-1} \boldsymbol{\mu}_{t-1} - \langle A > \Sigma_{t-1}^* \langle A^{\top} B > \mathbf{u}_t \rangle, \right) \tag{46}
$$

We approximate the forward message as a Gaussian in  $x_t$  using the first and second derivatives w.r.t.  $x_t$ 

$$
\alpha(\mathbf{x}_t) \quad \propto \quad \exp(<\log p(\mathbf{y}_t|\mathbf{x}_t) >_{q_\theta(\theta)}) \mathcal{N}(\mathbf{x}_t|\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\Sigma}}_t). \tag{47}
$$

where

$$
\langle \log p(\mathbf{y}_t | \mathbf{x}_t) \rangle_{q_\theta(\theta)} = \int \log p(\mathbf{y}_t | \mathbf{x}_t, \hat{C}, \mathbf{d}, \mathbf{h}^{(i)}) \mathcal{N}(\mathbf{h}^{(i)} | \boldsymbol{\mu}_{\mathbf{h}^{(i)}}, \Sigma_{\mathbf{h}^{(i)}}) d\mathbf{h}^{(i)},
$$
\n(48)

$$
= \int \left[ \mathbf{y}_t^\top (\hat{C}(\mathbf{x}_t + \mathbf{h}^{(i)}) + \mathbf{d}) - \mathbf{1}^\top \exp(\hat{C}(\mathbf{x}_t + \mathbf{h}^{(i)}) + \mathbf{d}) \right] \mathcal{N}(\mathbf{h}^{(i)} | \boldsymbol{\mu}_{\mathbf{h}^{(i)}}, \boldsymbol{\Sigma}_{\mathbf{h}^{(i)}}) d\mathbf{h}^{(i)}, \qquad (49)
$$

$$
= \sum_{s=1}^{p} \int \left[ (\mathbf{y}_t^\top \mathbf{e}_s)(\mathbf{x}_t^\top \hat{\mathbf{c}}_s) - \exp(\mathbf{x}_t^\top \hat{\mathbf{c}}_s + \mathbf{h}^{(i)\top} \hat{\mathbf{c}}_s + \mathbf{d}_s) \right] \mathcal{N}(\mathbf{h}^{(i)} | \boldsymbol{\mu}_{\mathbf{h}^{(i)}}, \Sigma_{\mathbf{h}^{(i)}}) d\mathbf{h}^{(i)}, \tag{50}
$$

$$
= \sum_{s=1}^{p} \left[ (\mathbf{y}_t^\top \mathbf{e}_s)(\mathbf{x}_t^\top \hat{\mathbf{c}}_s) - \exp(\mathbf{x}_t^\top \hat{\mathbf{c}}_s + \hat{\mathbf{c}}_s^\top \boldsymbol{\mu}_{\mathbf{h}^{(i)}} + \frac{1}{2} \hat{\mathbf{c}}_s^\top \boldsymbol{\Sigma}_{\mathbf{h}^{(i)}} \hat{\mathbf{c}}_s + \mathbf{d}_s) \right]
$$
(51)

The forward message at time  $t$  is approximately

$$
\alpha(\mathbf{x}_t) \approx \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t, \Sigma_t), \tag{52}
$$

$$
\mu_t = \tilde{\mu}_t + \tilde{\Sigma}_t \sum_{s=1}^P \left[ \mathbf{y}_t^T \mathbf{e}_s - \exp(\mathbf{x}_t^\top \hat{\mathbf{c}}_s + \hat{\mathbf{c}}_s^\top \boldsymbol{\mu}_{\mathbf{h}^{(i)}} + \frac{1}{2} \hat{\mathbf{c}}_s^\top \Sigma_{\mathbf{h}^{(i)}} \hat{\mathbf{c}}_s + \mathbf{d}_s) \right] \hat{\mathbf{c}}_s, \text{ where } \mathbf{x}_t = \boldsymbol{\mu}_t,
$$
\n(53)

$$
\Sigma_t^{-1} = \tilde{\Sigma}_t^{-1} + \sum_{s=1}^p \exp(\mathbf{x}_t^\top \hat{\mathbf{c}}_s + \hat{\mathbf{c}}_s^\top \boldsymbol{\mu}_{\mathbf{h}^{(i)}} + \frac{1}{2} \hat{\mathbf{c}}_s^\top \Sigma_{\mathbf{h}^{(i)}} \hat{\mathbf{c}}_s + \mathbf{d}_s) \hat{\mathbf{c}}_s \hat{\mathbf{c}}_s^\top, \text{ where } \mathbf{x}_t = \boldsymbol{\mu}_t. \tag{54}
$$

### Backward smoothing

We denote the backward message at each time  $t$  by

$$
\beta(\mathbf{x}_t) = p(\mathbf{y}_{t+1:T}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t|\boldsymbol{\eta}_t, \boldsymbol{\Psi}_t). \tag{55}
$$

We can obtain the recursion rules by considering  $\beta({\mathbf x}_{t-1})$ 

$$
\beta(\mathbf{x}_{t-1}) = \int d\mathbf{x}_t \beta(\mathbf{x}_t) \exp \left( \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_{\theta}(\theta)} \right),
$$
  
\n
$$
= \int d\mathbf{x}_t \exp \left( \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) \rangle_{q_{\theta}(\theta)} ) \left[ \beta(\mathbf{x}_t) \exp \left( \langle \log p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_{\theta}(\theta)} \right) \right],
$$
  
\n
$$
= \int d\mathbf{x}_t \exp \left( \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) \rangle_{q_{\theta}(\theta)} ) \mathcal{N}(\mathbf{x}_t|\tilde{\boldsymbol{\eta}}_t, \tilde{\boldsymbol{\Psi}}_t), \right)
$$
(56)

$$
\tilde{\eta}_t = \eta_t + \Psi_t \sum_{s=1}^p \left[ \mathbf{y}_t^T \mathbf{e}_s - \exp(\mathbf{x}_t^\top \hat{\mathbf{c}}_s + \hat{\mathbf{c}}_s^\top \boldsymbol{\mu}_{\mathbf{h}^{(i)}} + \frac{1}{2} \hat{\mathbf{c}}_s^\top \Sigma_{\mathbf{h}^{(i)}} \hat{\mathbf{c}}_s + \mathbf{d}_s) \right] \hat{\mathbf{c}}_s, \text{ where } \mathbf{x}_t = \tilde{\eta}_t,
$$
\n(57)

$$
\tilde{\Psi}_t^{-1} = \Psi_t^{-1} + \sum_{s=1}^p \exp(\mathbf{x}_t^\top \hat{\mathbf{c}}_s + \hat{\mathbf{c}}_s^\top \boldsymbol{\mu}_{\mathbf{h}^{(i)}} + \frac{1}{2} \hat{\mathbf{c}}_s^\top \Sigma_{\mathbf{h}^{(i)}} \hat{\mathbf{c}}_s + \mathbf{d}_s) \hat{\mathbf{c}}_s \hat{\mathbf{c}}_s^\top, \text{ where } \mathbf{x}_t = \tilde{\eta}_t. \tag{58}
$$

The first term in the integrand above is given by

$$
\langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) >_{q_{\theta}(\theta)} = -\frac{1}{2}(\mathbf{x}_t^\top \mathbf{x}_t - 2\mathbf{x}_t^\top (\langle A > \mathbf{x}_{t-1} + \langle B > \mathbf{u}_t \rangle)) - \frac{1}{2}(\mathbf{x}_{t-1}^\top \langle A^\top A > \mathbf{x}_{t-1} + 2\mathbf{x}_{t-1}^\top \langle A^\top B > \mathbf{u}_t \rangle) - \frac{1}{2} \mathbf{u}_t^\top \langle B^\top B > \mathbf{u}_t. \tag{59}
$$

Therefore, the integral is given by

$$
\int d\mathbf{x}_t \exp\left(\langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) >_{q_{\theta}(\theta)} \right) \mathcal{N}(\mathbf{x}_t|\tilde{\boldsymbol{\eta}}_t, \tilde{\boldsymbol{\Psi}}_t) = \tilde{Z} \int d\mathbf{x}_t \exp\left(-\frac{1}{2}\mathbf{x}_t^\top (I + \tilde{\boldsymbol{\Psi}}_t^{-1})\mathbf{x}_t + \mathbf{x}_t^\top (\langle A \rangle \mathbf{x}_{t-1} + \langle B \rangle \mathbf{u}_t + \tilde{\boldsymbol{\Psi}}_t^{-1} \tilde{\boldsymbol{\eta}}_t)\right)
$$

where (only showing the terms depending on  $\mathbf{x}_{t-1}$ )

$$
\tilde{Z} = -\frac{1}{2} (\mathbf{x}_{t-1}^\top < A^\top A > \mathbf{x}_{t-1} + 2\mathbf{x}_{t-1}^\top < A^\top B > \mathbf{u}_t) + \cdots \tag{60}
$$

After integrating out  $\mathbf{x}_t$  by formulating a Gaussian distribution  $\mathcal{N}(\mathbf{x}_t | \boldsymbol{\eta}_t^*, \Psi_t^*)$  where the mean and covariance are given by

$$
\Psi_t^{*-1} = I + \tilde{\Psi}_t^{-1},\tag{61}
$$

$$
\eta_t^* = \Psi_t^* (\mathbf{x}\_{t-1} + \mathbf{u}\_t + \tilde{\Psi}\_t^{-1} \tilde{\eta}\_t\), \tag{62}
$$

we obtain a quadratic function in  $\mathbf{x}_{t-1}$  (combining the remainder from the integral and  $\tilde{Z})$ 

$$
\frac{1}{2}(\mathbf{x}\_{t-1} + \mathbf{u}\_t + \tilde{\Psi}\_t^{-1} \tilde{\eta}\_t\)^{\top} \Psi\_t^\* \(\mathbf{x}\\_{t-1} + \mathbf{u}\\_t + \tilde{\Psi}\\_t^{-1} \tilde{\eta}\\_t\\) - \frac{1}{2} \\(\mathbf{x}\\_{t-1}^{\top} < A^{\top} A > \mathbf{x}\\_{t-1} + 2 \mathbf{x}\\_{t-1}^{\top} < A^{\top} B > \mathbf{u}\\_t\\)
$$
\n
$$
= -\frac{1}{2} (\mathbf{x}_{t-1}^{\top} ( - ^{\top} \Psi\_t^\* \\) \mathbf{x}\\_{t-1} - 2 \mathbf{x}\\_{t-1}^{\top} \\(^{\top} \Psi\\\_t^\\\* \mathbf{u}\\\_t + \tilde{\Psi}\\\_t^{-1} \tilde{\eta}\\\_t\\\) -  \mathbf{u}\\\_t\\\)\\\) + \cdots. \tag{63}
$$

Therefore, the backward message is approximately Gaussian with the mean and covariance given by

$$
\beta(\mathbf{x}_{t-1}) \approx \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\eta}_{t-1}, \Psi_{t-1}), \tag{64}
$$

$$
\Psi_{t-1}^{-1} = \langle A^{\top} A \rangle - \langle A \rangle^{\top} \Psi_t^* \langle A \rangle, \tag{65}
$$

$$
\eta_{t-1} = \Psi_{t-1}(^{\top} \Psi\_t^\* \( \mathbf{u}\_t + \tilde{\Psi}\_t^{-1} \tilde{\eta}\_t\) -  \mathbf{u}\_t\). \tag{66}
$$

#### Computing marginals of latent variables using  $\alpha$  and  $\beta$

Using the  $\alpha$  and  $\beta$  recursions in the forward/backward algorithm, we can compute the marginals of the latent variables.

$$
p(\mathbf{x}_t|\mathbf{y}_{1:T}) = p(\mathbf{x}_t|\mathbf{y}_{1:t}, \mathbf{y}_{t+1:T}),
$$
\n(67)

$$
\propto p(\mathbf{y}_{t+1:T}|\mathbf{x}_t,\mathbf{y}_{1:t})p(\mathbf{x}_t|\mathbf{y}_{1:t}) = p(\mathbf{y}_{t+1:T}|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t}) = \beta(\mathbf{x}_t)\alpha(\mathbf{x}_t),
$$
\n(68)

$$
\propto \mathcal{N}(\mathbf{x}_t|\boldsymbol{\omega}_t, \Upsilon_t) \tag{69}
$$

where

$$
\Upsilon_t^{-1} = \Psi_t^{-1} + \Sigma_t^{-1},\tag{70}
$$

$$
\omega_t = \Upsilon_t (\Psi_t^{-1} \eta_t + \Sigma_t^{-1} \mu_t). \tag{71}
$$

,

We also need to compute pairwise marginals of latent variables, given by

$$
p(\mathbf{x}_{t}, \mathbf{x}_{t+1} | \mathbf{y}_{1:T}) = p(\mathbf{x}_{t}, \mathbf{x}_{t+1} | \mathbf{y}_{1:t}, \mathbf{y}_{t+1}, \mathbf{y}_{t+2:T}),
$$
  
\n
$$
\propto p(\mathbf{y}_{t+1}, \mathbf{y}_{t+2:T} | \mathbf{x}_{t}, \mathbf{x}_{t+1}, \mathbf{y}_{1:t}) p(\mathbf{x}_{t+1} | \mathbf{x}_{t}, \mathbf{y}_{1:t}) p(\mathbf{x}_{t} | \mathbf{y}_{1:t}),
$$
  
\n
$$
\propto p(\mathbf{y}_{t+1} | \mathbf{x}_{t+1}) p(\mathbf{y}_{t+2:T} | \mathbf{x}_{t+1}) p(\mathbf{x}_{t+1} | \mathbf{x}_{t}) p(\mathbf{x}_{t} | \mathbf{y}_{1:t}),
$$
  
\n
$$
\propto \beta(\mathbf{x}_{t+1}) \exp \left( \langle \log(p(\mathbf{y}_{t+1} | \mathbf{x}_{t+1}) p(\mathbf{x}_{t+1} | \mathbf{x}_{t})) \rangle_{q_{\theta}(\theta)} \right) \alpha(\mathbf{x}_{t}),
$$
\n(72)

which are jointly Gaussian

$$
p\begin{pmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{pmatrix} = \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\omega}_t \\ \boldsymbol{\omega}_{t+1} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Upsilon}_t & \boldsymbol{\Upsilon}_{t,t+1} \\ \boldsymbol{\Upsilon}_{t,t+1}^T & \boldsymbol{\Upsilon}_{t+1} \end{bmatrix} \right).
$$
(73)

To compute the cross-covariance  $\Upsilon_{t,t+1}$ , we first compute the second derivatives w.r.t.  $[\mathbf{x}_t \; \mathbf{x}_{t+1}]^T$ :

$$
\frac{\partial^2 \log \int d\theta q_{\theta}(\theta) p(\mathbf{x}_t, \mathbf{x}_{t+1} | \mathbf{y}_{1:T})}{\partial [\mathbf{x}_t | \mathbf{x}_{t+1}]^2} = -\begin{bmatrix} \Sigma_t^{*-1} & -\langle A \rangle^T\\ -\langle A \rangle & \Psi_{t+1}^{-1} + I + W_{t+1} \end{bmatrix},\tag{74}
$$

where

$$
W_{t+1} = -\frac{\partial^2}{\partial \mathbf{x}_{t+1} \mathbf{x}_{t+1}^\top} < \log p(\mathbf{y}_{t+1} | \mathbf{x}_{t+1}) >_{q(\theta)},\tag{75}
$$

$$
= \sum_{s=1}^{p} \exp(\mathbf{x}_{t+1}^\top \hat{\mathbf{c}}_s + \hat{\mathbf{c}}_s^\top \boldsymbol{\mu}_{\mathbf{h}^{(i)}} + \frac{1}{2} \hat{\mathbf{c}}_s^\top \boldsymbol{\Sigma}_{\mathbf{h}^{(i)}} \hat{\mathbf{c}}_s + \mathbf{d}_s) \hat{\mathbf{c}}_s \hat{\mathbf{c}}_s^\top
$$
\n(76)

evaluated at  $\mathbf{x}_{t+1} = \boldsymbol{\omega}_{t+1}$ . By negating and inverting the matrix, and using the *Schur* complement, we can obtain  $\Upsilon_{t,t+1}$ ,

$$
\Upsilon_{t,t+1} = -(\Sigma_t^{*-1} - \langle A \rangle^T (\Psi_{t+1}^{-1} + I + W_{t+1})^{-1} \langle A \rangle)^{-1} (-\langle A \rangle^T) (\Psi_{t+1}^{-1} + I + W_{t+1})^{-1}.
$$
\n(77)

### Computing sufficient statistics of latent variables

Using  $q_{\mathbf{x}}(\mathbf{x}_{0:T})$ , we can compute the sufficient statistics of latent variables (that are used in M step).

$$
W_A = \sum_{t=1}^T <\mathbf{x}_{t-1}\mathbf{x}_{t-1}^T> = \sum_{t=1}^T \Upsilon_{t-1} + \omega_{t-1}\omega_{t-1}^T, \qquad S_A = \sum_{t=1}^T <\mathbf{x}_{t-1}\mathbf{x}_t^T> = \sum_{t=1}^T \Upsilon_{t-1,t} + \omega_{t-1}\omega_t^T,
$$
(78)

$$
W_C = \sum_{t=1}^T <\mathbf{x}_t \mathbf{x}_t^T> = \sum_{t=1}^T \Upsilon_t + \omega_t \omega_t^T, \qquad S_C = \sum_{t=1}^T <\mathbf{x}_t > \mathbf{y}_t^T = \sum_{t=1}^T \omega_t \mathbf{y}_t^T.
$$
 (79)

## Hyperaparameter estimation

We take the derivatives of the variational lower bound w.r.t. each hyperparameter to obtain update rules. The lower bound is simplified as below:

$$
\log p(\mathbf{y}_{1:T}^{(1:r)}) \geq \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} \, q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log \frac{p(\theta, \mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)})}{q(\theta, \mathbf{x}_{0:T}^{(1:r)})},
$$
\n
$$
= \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} \, q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log p(\mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)}) \theta) - \int d\mathbf{x}_{0:T}^{(1:r)} \, q(\mathbf{x}_{0:T}^{(1:r)}) \log q(\mathbf{x}_{0:T}^{(1:r)}) + \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} \, q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log \frac{p(\theta)}{q(\theta)},
$$
\n
$$
= \sum_{i=1}^{r} \log Z'_{(i)} + \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} \, q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log \frac{p(\theta)}{q(\theta)},
$$
\n(80)

where the last line follows from the equality

$$
-\int d\mathbf{x}_{0:T}^{(1:r)} q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)}) \log q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)}) = -\int d\mathbf{x}_{0:T}^{(1:r)} q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)}) \mathbb{E}_{q_{\theta}(\theta)} \log p(\mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)} | \theta) + \sum_{i=1}^{r} \log Z'_{(i)}.
$$
 (81)

So, we need to consider the second term in RHS of the lower bound for hyperparameter update (the integration w.r.t. x is omitted, since the integrand is independent of x)

$$
\int d\mathbf{a} \, d\mathbf{b} \, d\mathbf{h}^{(1:r)} q(\mathbf{a}, \mathbf{b}) q(\mathbf{h}^{(1:r)}) \, \log \frac{p(\mathbf{a}, \mathbf{b}, \mathbf{h}^{(1:r)})}{q(\mathbf{a}, \mathbf{b}, \mathbf{h}^{(1:r)})} \quad = \quad -KL(\mathbf{a}, \mathbf{b}) - KL(\mathbf{h}^{(1:r)}), \tag{82}
$$

where the first term on RHS is given by

$$
KL(\mathbf{a}, \mathbf{b}) = \int d\mathbf{a} \, d\mathbf{b} q(\mathbf{a}, \mathbf{b}) \, \log \frac{q(\mathbf{a}, \mathbf{b})}{p(\mathbf{a}, \mathbf{b})},\tag{83}
$$

$$
= \int d\mathbf{a} \, d\mathbf{b} \mathcal{N}(\boldsymbol{\mu}_{\mathbf{a},\mathbf{b}}, \Lambda_{\mathbf{a},\mathbf{b}}) \, \log \frac{\mathcal{N}(\boldsymbol{\mu}_{\mathbf{a},\mathbf{b}}, \Lambda_{\mathbf{a},\mathbf{b}})}{\mathcal{N}(0, \tilde{\Lambda}_{\mathbf{a},\mathbf{b}})},
$$
(84)

$$
= -\frac{1}{2}\log|\tilde{\Lambda}_{\mathbf{a},\mathbf{b}}^{-1}\Lambda_{\mathbf{a},\mathbf{b}}| + \frac{1}{2}\text{Tr}[\tilde{\Lambda}_{\mathbf{a},\mathbf{b}}^{-1}(\Lambda_{\mathbf{a},\mathbf{b}} + \mu_{\mathbf{a},\mathbf{b}}\mu_{\mathbf{a},\mathbf{b}}^{\top})] \tag{85}
$$

where the prior covariance on  $(a, b)$  is denoted by  $\tilde{\Lambda}_{a,b} = [\alpha^{-1} I \ 0; 0 \ \beta^{-1} I]$  and the posterior mean and covariance on  $(a, b)$ are denoted by  $\mu_{a,b} = [\mu_a; \mu_b]$  and  $\Lambda_{a,b} = [\Sigma_a \Sigma_{a,b}; \Sigma_{a,b}^\top \Sigma_b]$ , respectively. We minimise the KL divergence for updating α, β.

The second term on RHS is given by

$$
KL(\mathbf{h}^{(1:r)}) = \int d\mathbf{h}^{(1:r)} q_{\mathbf{h}}(\mathbf{h}^{(1:r)}) \log \frac{q_{\mathbf{h}}(\mathbf{h}^{(1:r)})}{p(\mathbf{h}^{(1:r)}|\sigma^2, \tau^2)},
$$
\n(86)

$$
= \int d\mathbf{h}^{(1:r)} \mathcal{N}(\mathbf{h}^{(1:r)}|\boldsymbol{\mu_h}, \boldsymbol{\Sigma_h}) \log \frac{\mathcal{N}(\mathbf{h}^{(1:r)}|\boldsymbol{\mu_h}, \boldsymbol{\Sigma_h})}{p(\mathbf{h}^{(1:r)}|\mathbf{m_h}, K)},
$$
\n(87)

$$
= -\frac{1}{2}\log|K^{-1}\Sigma_{\mathbf{h}}| + \frac{1}{2}\text{Tr}\left[K^{-1}(\Sigma_{\mathbf{h}} + (\boldsymbol{\mu}_{\mathbf{h}} - \mathbf{m}_{\mathbf{h}})(\boldsymbol{\mu}_{\mathbf{h}} - \mathbf{m}_{\mathbf{h}})^{\top})\right] + const.
$$
\n(88)

The first derivative w.r.t. kernel parameters (denoted by  $\alpha = \{\sigma^2, \tau^2\}$ ) is given by

$$
\frac{\partial}{\partial \alpha} KL(\mathbf{h}^{(1:r)}) = \frac{1}{2} \text{Tr} \left( K^{-1} \frac{\partial K}{\partial \alpha} \right) - \frac{1}{2} \text{Tr} \left( K^{-1} \frac{\partial K}{\partial \alpha} K^{-1} (\Sigma_{\mathbf{h}} + (\boldsymbol{\mu}_{\mathbf{h}} - \mathbf{m}_{\mathbf{h}}) (\boldsymbol{\mu}_{\mathbf{h}} - \mathbf{m}_{\mathbf{h}})^{\top}) \right), \tag{89}
$$

$$
= \frac{1}{2} \text{Tr} \left( K^{-1} \frac{\partial K}{\partial \alpha} (I - K^{-1} (\Sigma_{\mathbf{h}} + (\mu_{\mathbf{h}} - \mathbf{m}_{\mathbf{h}}) (\mu_{\mathbf{h}} - \mathbf{m}_{\mathbf{h}})^{\top})) \right), \tag{90}
$$

where the first derivative of  $K(i, j)$  w.r.t.  $\alpha$  is given by

$$
\frac{\partial}{\partial \tau^2} K(i,j) = \frac{1}{2\tau^4} (i-j)^2 (\sigma^2 + \epsilon \delta_{ij}) \exp\left(-\frac{1}{2\tau^2} (i-j)^2\right) I_{k^2} = \frac{1}{2\tau^4} (i-j)^2 K(i,j),\tag{91}
$$

$$
\frac{\partial}{\partial \sigma^2} K(i,j) = \exp\left(-\frac{1}{2\tau^2}(i-j)^2\right) I_{k^2}.
$$
\n(92)

We update  $\alpha$  numerically using the derivative expression above.

# Model II: nonstationarity in neural dynamics

# Basic setup

Likelihood:  $\mathbf{y}_t \in \mathbb{R}^p$ ,  $\mathbf{x}_t \in \mathbb{R}^k$ ,  $C \in \mathbb{R}^{p \times k}$ 

$$
p(\mathbf{y}_t|\mathbf{x}_t, C, \mathbf{d}) = \text{Poiss}(\mathbf{y}_t|\exp(C\mathbf{x}_t + \mathbf{d})).
$$

Latent dynamics:  $A \in \mathbb{R}^{k \times k}$ 

$$
p(\mathbf{x}_t|\mathbf{x}_{t-1}, A) = \mathcal{N}(\mathbf{x}_t|A\mathbf{x}_{t-1}, I).
$$

Parameters in this model:  $\Theta = \{A, C\}$ . For simplicity, we will fix **d** to its maximum likelihood estimate. Vectorized notations:  $\mathbf{a} = \text{vec}(A^{\top}) \in \mathbb{R}^{k^2}$  and  $\mathbf{c} = \text{vec}(C^{\top}) \in \mathbb{R}^{pk}$ .

#### Priors:

$$
p(\mathbf{c}|\gamma) = \mathcal{N}(\mathbf{c}|\mathbf{0}, \gamma^{-1}\mathbf{I})
$$
\n(93)

Assuming a to be temporally evolving across recordings where the recording index is  $i = \{1, \dots, r\}$ :

$$
\mathbf{a}^{(i)} \sim \mathcal{GP}(\bar{\mathbf{a}}, K(i,j)) \tag{94}
$$

where we denote the (vector) mean and (matrix) covariance functions by  $\bar{a}$  and  $K(i, j)$ , respectively, where the  $(i, j)$ th block of the covariance matrix is given by

$$
K(i,j) = (\sigma^2 + \epsilon \delta_{i,j}) \exp\left(-\frac{1}{2\tau^2}(i-j)^2\right) I_{k^2}.
$$
\n
$$
(95)
$$

The hyperparameters in total are  $\Phi = {\bar{\mathbf{a}}, \sigma^2, \tau^2, \gamma}.$ 

# Variational lower bound

The marginal likelihood of the observations is lower bounded by

$$
\log p(\mathbf{y}_{1:T}^{(1:r)}) \geq \int d\theta \ d\mathbf{x}_{0:T}^{(1:r)} \ q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \ \log \frac{p(\theta, \mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)})}{q(\theta, \mathbf{x}_{0:T}^{(1:r)})},
$$
\n(96)

where the approximate posterior factories

$$
q(\theta, \mathbf{x}_{0:T}^{(1:r)}) = q_{\theta}(\theta) \prod_{i=1}^{r} q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)}),
$$
\n(97)

and we assume  $q_{\theta}(\theta) = q_{\mathbf{a}}(\mathbf{a}^{(1:r)})q_{\mathbf{c}}(\mathbf{c}).$ 

# Bayesian Laplace propagation

## Posterior over parameters

We compute  $q_{\theta}(\theta)$  by integrating out latent variables from the total log joint distribution:

$$
\log q_{\theta}(\theta) = \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)})}\left[\log p(\mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)}, \theta)\right] + const,
$$
\n
$$
= \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)})}\left[\log p(\mathbf{y}_{1:T}^{(1:r)} | \mathbf{x}_{0:T}^{(1:r)}, \theta) + \log p(\mathbf{x}_{0:T}^{(1:r)} | \theta) + \log p(\theta)\right] + const,
$$
\n
$$
= \sum_{i=1}^{r} \left[ \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)})} (\sum_{t=1}^{T} (\log p(\mathbf{y}_{t}^{(i)} | \mathbf{x}_{t}^{(i)}, \mathbf{c}) + \log p(\mathbf{x}_{t}^{(i)} | \mathbf{x}_{t-1}^{(i)}, \mathbf{a}^{(i)})) \right] + \log p(\mathbf{a}^{(1:r)} | \mathbf{\bar{a}}^{(1:r)}, K) + \log p(\mathbf{c} | \gamma) + const,
$$
\n(98)

where  $\bar{\mathbf{a}}^{(1:r)}$  is a vector of r repeating  $\bar{\mathbf{a}}$ .

# 1. approximate posterior over  $a^{(1:r)}$

$$
\log q_{\mathbf{a}}(\mathbf{a}^{(1:r)}) = \sum_{i=1}^{r} \left[ \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)})} \sum_{t=1}^{T} \log p(\mathbf{x}_{t}^{(i)} | \mathbf{x}_{t-1}^{(i)}, \mathbf{a}^{(i)}) \right] + \log p(\mathbf{a}^{(1:r)} | \mathbf{\bar{a}}^{(1:r)}, K) + const,
$$
\n(99)

$$
= -\frac{1}{2} (\mathbf{a}^{(1:r)} \top H \mathbf{a}^{(1:r)} - 2 \mathbf{a}^{(1:r)} \top \mathbf{s}) - \frac{1}{2} (\mathbf{a}^{(1:r)} - \bar{\mathbf{a}}^{(1:r)})^T K^{-1} (\mathbf{a}^{(1:r)} - \bar{\mathbf{a}}^{(1:r)})
$$
(100)

where the matrix  $H$  and the vector  $s$  are given by

$$
H_{\mathbf{a}} = \begin{pmatrix} W_{A^{(1)}}^{bd}, & 0, & 0, & \dots & 0 \\ 0, & W_{A^{(2)}}^{bd}, & 0, & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0, & \dots & \dots, & 0, & W_{A^{(r)}}^{bd} \end{pmatrix}, \quad \mathbf{s} = \begin{pmatrix} \text{vec}(S_{A^{(1)}}) \\ \vdots \\ \text{vec}(S_{A^{(r)}}) \end{pmatrix}
$$
(101)

and  $W_{A^{(i)}}^{bd} = I_k \otimes W_{A^{(i)}},$  where  $W_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \mathbf{x}_{t-1}^{(i)} |^{\top} \rangle$ , and  $S_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \mathbf{x}_{t}^{(i)} |^{\top} \rangle$ . Therefore, the approximate posterior over  $\mathbf{a}^{(1:r)}$  is given by

$$
q(\mathbf{a}^{(1:r)}) = \mathcal{N}(\boldsymbol{\mu}_\mathbf{a}, \boldsymbol{\Sigma}_\mathbf{a}), \tag{102}
$$

$$
\Sigma_{\mathbf{a}}^{-1} = K^{-1} + H_{\mathbf{a}},\tag{103}
$$

$$
\mu_{\mathbf{a}} = \Sigma_{\mathbf{a}} (K^{-1} \bar{\mathbf{a}}^{(1:r)} + \mathbf{s}). \tag{104}
$$

 $\text{So,} < A^{(i)}>=\left[\text{reshape}(\mu_{\mathbf{a}}((i-1)k^2+1:ik^2), k, k)\right]^\top$  and  $\Sigma_{A^{(i)}}$  is the first  $k\times k$  matrix of  $\Sigma_{\mathbf{a}}((i-1)k^2+1:ik^2,(i-1)k^2+1:$  $ik<sup>2</sup>$ ). In addition to the mean and covariance of  $A<sup>(i)</sup>$ , we also need the following quantity in VBE step:

$$
\langle A^{(i)\top} A^{(i)} \rangle = \langle A^{(i)} \rangle^{\top} \langle A^{(i)} \rangle + k \Sigma_{A^{(i)}}. \tag{105}
$$

### 12

## 2. Computing  $q_c(c)$

Similarly, we write down all the terms in  $\log p(\mathbf{x}_{0:T}, \mathbf{y}_{1:T}, \theta)$  that depend on **c**:

$$
\log q_{\mathbf{c}}(\mathbf{c}) = \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \left[ \sum_{t=1}^{T} \log p(\mathbf{y}_{t}^{(i)} | \mathbf{x}_{t}^{(i)}, \mathbf{c}) \right] + \log p(\mathbf{c}|\gamma) + const,
$$
  
\n
$$
= \sum_{i=1}^{r} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \left[ \sum_{t=1}^{T} (\mathbf{y}_{t}^{(i)\top} (C\mathbf{x}_{t}^{(i)} + \mathbf{d}) - \mathbf{1}^{\top} \exp(C\mathbf{x}_{t}^{(i)} + \mathbf{d})) \right] - \frac{1}{2} \gamma \mathbf{c}^{T} \mathbf{c} + const,
$$
  
\n
$$
= \sum_{i=1}^{r} \left[ \mathbf{c}^{\top} \text{vec}(S_{C^{(i)}}) - \sum_{s=1}^{p} \mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})} \{ \sum_{t=1}^{T} \exp(c_{s}^{\top} \mathbf{x}_{t}^{(i)} + \mathbf{d}_{s}) \} \right] - \frac{1}{2} \gamma \mathbf{c}^{\top} \mathbf{c} + const,
$$
(106)

where each row of  $C$  is denoted by  $\mathbf{c}_s$  and the sufficient statistic is denoted by

$$
S_{C^{(i)}} = \sum_{t=1}^{T} < \mathbf{x}_t^{(i)} > \mathbf{y}_t^{(i)\top} \tag{107}
$$

Assuming the approximate posterior over latent variables is multivariate Gaussian with marginals  $q(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t | \boldsymbol{\omega}_t, \Upsilon_t)$ , the expectation of the exponential term in eq. 106 is given by

$$
\mathbb{E}_{q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)})}\left\{\sum_{t=1}^{T} \exp(\mathbf{c}_{s}^{\top}\mathbf{x}_{t}^{(i)})\right\} = \int d\mathbf{x}_{1:T}^{(i)} q_{\mathbf{x}}(\mathbf{x}_{1:T}^{(i)}) \exp(\mathbf{c}_{s}^{\top}\mathbf{x}_{1}^{(i)} + \cdots + \mathbf{c}_{s}^{\top}\mathbf{x}_{T}^{(i)}),
$$
\n
$$
= \sum_{t=1}^{T} \exp(\mathbf{c}_{s}^{\top}\boldsymbol{\omega}_{t}^{(i)} + \frac{1}{2}\mathbf{c}_{s}^{\top}\boldsymbol{\Upsilon}_{t}^{(i)}\mathbf{c}_{s}). \tag{108}
$$

Therefore, the log joint distribution is given by

$$
\log q_{\mathbf{c}}(\mathbf{c}) = \sum_{i=1}^{r} \left[ \mathbf{c}^{\top} \text{vec}(S_{C^{(i)}}) - \sum_{s=1}^{p} \sum_{t=1}^{T} \exp(\mathbf{c}_{s}^{\top} \boldsymbol{\omega}_{t}^{(i)} + \frac{1}{2} \mathbf{c}_{s}^{\top} \boldsymbol{\Upsilon}_{t}^{(i)} \mathbf{c}_{s} + \mathbf{d}_{s}) \right] - \frac{1}{2} \gamma \mathbf{c}^{\top} \mathbf{c} + const.
$$
 (109)

We approximate  $q_C(C)$  to a Gaussian distribution from the first/second derivatives of eq. 109 w.r.t.  $\mathbf{c}_s$ ,

$$
q_C(C) = \prod_{s=1}^p \mathcal{N}(\mathbf{c}_s | \boldsymbol{\mu}_{\mathbf{c}_s}, \Sigma_{\mathbf{c}_s})
$$
\n(110)

$$
\boldsymbol{\mu}_{\mathbf{c}_s} = \frac{1}{\gamma} \sum_{i=1}^r \left[ S_{C^{(i)}} \mathbf{e}_s - \sum_{t=1}^T [\boldsymbol{\omega}_t^{(i)} + \boldsymbol{\Upsilon}_t^{(i)} \mathbf{c}_s] \exp(\mathbf{c}_s^T \boldsymbol{\omega}_t^{(i)} + \frac{1}{2} \mathbf{c}_s^T \boldsymbol{\Upsilon}_t^{(i)} \mathbf{c}_s + \mathbf{d}_s) \right], \text{ where } \mathbf{c}_s = \boldsymbol{\mu}_{\mathbf{c}_s},
$$
\n(111)

$$
\Sigma_{\mathbf{c}_s}^{-1} = \gamma I + \sum_{i=1}^r \sum_{t=1}^T (\Upsilon_t^{(i)} + (\omega_t^{(i)} + \Upsilon_t^{(i)} \mathbf{c}_s)(\omega_t^{(i)} + \Upsilon_t^{(i)} \mathbf{c}_s)^\top) \exp(\mathbf{c}_s^T \omega_t^{(i)} + \frac{1}{2} \mathbf{c}_s^T \Upsilon_t^{(i)} \mathbf{c}_s + \mathbf{d}_s), \text{ where } \mathbf{c}_s = \boldsymbol{\mu}_{\mathbf{c}_s} \cdot (112)
$$

#### 3. ML estimate of d

The ML estimate of **d** given the mean of C (denoted by  $\hat{C}$ ) is closed form:

$$
\hat{\mathbf{d}} = \arg \max_{\mathbf{d}} \sum_{i=1}^{r} \left[ \hat{C} S_{C^{(i)}} + \sum_{t=1}^{T} \left( \mathbf{y}_{t}^{(i) \top} \mathbf{d} - \mathbf{1}^{\top} \exp(\hat{C} \boldsymbol{\omega}_{t}^{(i)} + \frac{1}{2} \text{diag}(\hat{C} \Upsilon_{t}^{(i)} \hat{C}^{\top}) + \mathbf{d}) \right) \right],
$$
  
\n
$$
= \log(\sum_{i=1}^{r} \sum_{t=1}^{T} \mathbf{y}_{t}^{(i)}) - \log(\sum_{i=1}^{r} \sum_{t=1}^{T} \exp(\hat{C} \boldsymbol{\omega}_{t}^{(i)} + \frac{1}{2} \text{diag}(\hat{C} \Upsilon_{t}^{(i)} \hat{C}^{\top})). \tag{113}
$$

## Posterior over latent variables

We compute  $q_{\mathbf{x}}(\mathbf{x}_{0:qT})$  by

$$
\sum_{i=1}^{r} \log q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(i)}) = \sum_{i=1}^{r} \mathbb{E}_{q_{\theta}(\theta)} \log p(\theta, \mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)}) + const,\n= \sum_{i=1}^{r} \mathbb{E}_{q_{\theta}(\theta)} \log p(\mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)} | \theta) - \sum_{i=1}^{r} \log Z'_{(i)},
$$
\n(114)

where the normalization constant is given by

$$
Z'_{(i)} = \int d\mathbf{x}_{0:T}^{(i)} \exp\left(\mathbb{E}_{q_{\theta}(\theta)} \log p(\mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)} | \theta)\right).
$$
\n(115)

The complete-data log likelihood in the ith recording is written as

$$
\log p(\mathbf{x}_{0:T}^{(i)}, \mathbf{y}_{1:T}^{(i)} | \theta) = \sum_{t=1}^{T} \{ \log p(\mathbf{y}_t^{(i)} | \mathbf{x}_t^{(i)}, C, \mathbf{d}) + \log p(\mathbf{x}_t^{(i)} | \mathbf{x}_{t-1}^{(i)}, A^{(i)}) \},
$$
(116)

which tells us that the log posterior over latent variables is quadratic in each  $x_t$ . This enables us to use the sequential update of the posterior over latent variables. We will also use the following sequential forward/backward algorithm for each recording in parallel. In the following, the recording index i is removed for notational cleanness.

### Forward filtering

We denote the posterior over the latent variables at each time  $t$  by

$$
\alpha(\mathbf{x}_t) \propto \int d\mathbf{x}_{t-1} \alpha(\mathbf{x}_{t-1}) \exp \left[ \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_\theta(\theta)} \right], \tag{117}
$$

$$
\propto \exp(\langle \log p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_\theta(\theta)} \left\{ \int d\mathbf{x}_{t-1} \alpha(\mathbf{x}_{t-1}) \exp\left(\langle \log (p(\mathbf{x}_t|\mathbf{x}_{t-1}) \rangle_{q(\theta)}) \right) \right\}.
$$
 (118)

Assuming  $\alpha(\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{t-1}, \Sigma_{t-1})$ , the integral is analytically tractable since the second part in the integrand is also quadratic in  $\mathbf{x}_{t-1}$ :

$$
\exp[-\tfrac{1}{2}(\mathbf{x}_{t-1}^\top < A^\top A > \mathbf{x}_{t-1} - 2\mathbf{x}_{t-1}^\top < A > \top \mathbf{x}_t + \mathbf{x}_t^\top \mathbf{x}_t)].
$$

The integrand is summarised as

$$
\alpha(\mathbf{x}_{t-1}) \exp\left(\langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) >_{q(\theta)} \right) = Z\mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\mu}_{t-1}^*, \boldsymbol{\Sigma}_{t-1}^*), \tag{119}
$$

$$
\Sigma_{t-1}^{*-1} = \Sigma_{t-1}^{-1} + \langle A^{\top} A \rangle, \tag{120}
$$

$$
\mu_{t-1}^* = \Sigma_{t-1}^*(\Sigma_{t-1}^{-1}\mu_{t-1} + \langle A \rangle^\top \mathbf{x}_t), \tag{121}
$$

and the remaining term  $Z$  is given by:

$$
Z = \exp[\frac{1}{2}\mu_{t-1}^* \Gamma_{t-1}^{*-1} \mu_{t-1}^*], \qquad (122)
$$

where

$$
\frac{1}{2}\mu_{t-1}^* \mathbb{E}_{t-1}^{* - 1} \mu_{t-1}^* = \frac{1}{2} (\Sigma_{t-1}^{-1} \mu_{t-1} + \langle A \rangle^\top \mathbf{x}_t)^\top \Sigma_{t-1}^* (\Sigma_{t-1}^{-1} \mu_{t-1} + \langle A \rangle^\top \mathbf{x}_t),
$$
  
\n
$$
= \frac{1}{2} (\mathbf{x}_t^\top \langle A \rangle \Sigma_{t-1}^* \langle A \rangle^\top \mathbf{x}_t + 2 \mathbf{x}_t^\top \langle A \rangle \Sigma_{t-1}^* \Sigma_{t-1}^{-1} \mu_{t-1}) + \cdots.
$$

Therefore,  $Z$  is proportional to a Gaussian in  $\mathbf{x}_t$  :

$$
Z \propto \mathcal{N}(\mathbf{x}_t | \tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\Sigma}}_t), \tag{123}
$$

$$
\tilde{\Sigma}_t^{-1} = I - \langle A \rangle \Sigma_{t-1}^* \langle A \rangle^\top, \tag{124}
$$

$$
\tilde{\mu}_t = \tilde{\Sigma}_t < A > \Sigma_{t-1}^* \Sigma_{t-1}^{-1} \mu_{t-1},\tag{125}
$$

We approximate the forward message as a Gaussian in  $x_t$  using the first and second derivatives w.r.t.  $x_t$ 

 $\alpha(\mathbf{x}_t) \quad \propto \quad \exp(<\log p(\mathbf{y}_t|\mathbf{x}_t) >_{q_\theta(\theta)}) \mathcal{N}(\mathbf{x}_t|\tilde{\boldsymbol{\mu}}_t, \tilde{\boldsymbol{\Sigma}}_t).$  $(126)$ 

The forward message at time  $t$  is approximately

$$
\alpha(\mathbf{x}_t) \approx \mathcal{N}(\mathbf{x}_t | \boldsymbol{\mu}_t, \Sigma_t),
$$
\n(127)

$$
\mu_t = \tilde{\mu}_t + \tilde{\Sigma}_t \sum_{s=1}^P \left[ (\mathbf{y}_t^T \mathbf{e}_s) \mu_{\mathbf{c}_s} - (\mu_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_t) e^{\mathbf{x}_t^T \mu_{\mathbf{c}_s} + \frac{1}{2} \mathbf{x}_t^T \Sigma_{\mathbf{c}_s} \mathbf{x}_t + \mathbf{d}_s} \right], \text{ where } \mathbf{x}_t = \mu_t,
$$
\n(128)

$$
\Sigma_t^{-1} = \tilde{\Sigma}_t^{-1} + \sum_{s=1}^p \left[ \Sigma_{\mathbf{c}_s} + (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_t) (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_t)^T \right] e^{\mathbf{x}_t^T \boldsymbol{\mu}_{\mathbf{c}_s} + \frac{1}{2} \mathbf{x}_t^T \Sigma_{\mathbf{c}_s} \mathbf{x}_t + \mathbf{d}_s}, \text{ where } \mathbf{x}_t = \boldsymbol{\mu}_t. \tag{129}
$$

#### Backward smoothing

We denote the backward message at each time  $t$  by

$$
\beta(\mathbf{x}_t) = p(\mathbf{y}_{t+1:T}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t|\boldsymbol{\eta}_t, \boldsymbol{\Psi}_t). \tag{130}
$$

We can obtain the recursion rules by considering  $\beta(\mathbf{x}_{t-1})$ 

$$
\beta(\mathbf{x}_{t-1}) = \int d\mathbf{x}_t \beta(\mathbf{x}_t) \exp \left( \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1})p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_{\theta}(\theta)} \right),
$$
  
\n
$$
= \int d\mathbf{x}_t \exp \left( \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) \rangle_{q_{\theta}(\theta)} ) [\beta(\mathbf{x}_t) \exp \left( \langle \log p(\mathbf{y}_t|\mathbf{x}_t)) \rangle_{q_{\theta}(\theta)} \right) ] ,
$$
  
\n
$$
= \int d\mathbf{x}_t \exp \left( \langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) \rangle_{q_{\theta}(\theta)} ) \mathcal{N}(\mathbf{x}_t|\tilde{\boldsymbol{\eta}}_t, \tilde{\boldsymbol{\Psi}}_t), \right. (131)
$$

assuming  $\beta(\mathbf{x}_T) = 1$ . The Gaussian  $p(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t | \tilde{\boldsymbol{\eta}}_t, \tilde{\boldsymbol{\Psi}}_t)$  is obtained by computing the first and second derivatives w.r.t.  $\mathbf{x}_t,$ 

$$
\tilde{\eta}_t = \eta_t + \Psi_t \sum_{s=1}^p \left[ (\mathbf{y}_t^T \mathbf{e}_s) \boldsymbol{\mu}_{\mathbf{c}_s} - (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_t) e^{\mathbf{x}_t^T \boldsymbol{\mu}_{\mathbf{c}_s} + \frac{1}{2} \mathbf{x}_t^T \Sigma_{\mathbf{c}_s} \mathbf{x}_t + \mathbf{d}_s} \right], \text{ where } \mathbf{x}_t = \tilde{\eta}_t,
$$
\n(132)

$$
\tilde{\Psi}_t^{-1} = \Psi_t^{-1} + \sum_{s=1}^p \left[ \Sigma_{\mathbf{c}_s} + (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_t) (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_t)^T \right] e^{\mathbf{x}_t^T \boldsymbol{\mu}_{\mathbf{c}_s} + \frac{1}{2} \mathbf{x}_t^T \Sigma_{\mathbf{c}_s} \mathbf{x}_t + \mathbf{d}_s}, \quad \text{where } \mathbf{x}_t = \tilde{\boldsymbol{\eta}}_t.
$$
\n(133)

The first term in the integrand in eq. 131 is given by

$$
\langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) >_{q_\theta(\theta)} = -\frac{1}{2}(\mathbf{x}_t^T \mathbf{x}_t - 2\mathbf{x}_t^T \langle A > \mathbf{x}_{t-1}) - \frac{1}{2}\mathbf{x}_{t-1}^T \langle A^T A > \mathbf{x}_{t-1}.\tag{134}
$$

Therefore, the integral is given by

$$
\int d\mathbf{x}_t \exp\left(\langle \log(p(\mathbf{x}_t|\mathbf{x}_{t-1}) >_{q_{\theta}(\theta)} \right) \mathcal{N}(\mathbf{x}_t|\tilde{\boldsymbol{\eta}}_t, \tilde{\boldsymbol{\Psi}}_t) = \tilde{Z} \int d\mathbf{x}_t \exp\left(-\frac{1}{2}\mathbf{x}_t^T(I + \tilde{\boldsymbol{\Psi}}_t^{-1})\mathbf{x}_t + \mathbf{x}_t^T(\langle A \rangle \mathbf{x}_{t-1} + \tilde{\boldsymbol{\Psi}}_t^{-1} \tilde{\boldsymbol{\eta}}_t)\right)
$$

where (only showing the terms depending on  $\mathbf{x}_{t-1}$ )

$$
\tilde{Z} = -\frac{1}{2}\mathbf{x}_{t-1}^T < A^T A > \mathbf{x}_{t-1} + \cdots
$$
\n(135)

After integrating out  $\mathbf{x}_t$  by formulating a Gaussian distribution  $\mathcal{N}(\mathbf{x}_t | \boldsymbol{\eta}_t^*, \Psi_t^*)$  where the mean and covariance are given by

$$
\Psi_t^{*-1} = I + \tilde{\Psi}_t^{-1},\tag{136}
$$

$$
\boldsymbol{\eta}_t^* = \Psi_t^* \langle \langle A \rangle \mathbf{x}_{t-1} + \tilde{\Psi}_t^{-1} \tilde{\boldsymbol{\eta}}_t \rangle, \tag{137}
$$

we obtain a quadratic function in  $\mathbf{x}_{t-1}$  (combining the remainder from the integral and  $\tilde{Z})$ 

$$
\frac{1}{2}(\mathbf{x}\_{t-1}+\tilde{\Psi}\_{t}^{-1}\tilde{\pmb{\eta}}\_{t}\)^{T}\Psi\_{t}^{\*}\(\mathbf{x}\\_{t-1}+\tilde{\Psi}\\_{t}^{-1}\tilde{\pmb{\eta}}\\_{t}\\)-\frac{1}{2}\mathbf{x}\\_{t-1}^{T}\mathbf{x}\\_{t-1}
$$
\n
$$
= -\frac{1}{2}(\mathbf{x}_{t-1}^{T}(-)\mathbf{x}_{t-1}-2\mathbf{x}_{t-1}^{T}^{T}\Psi\_{t}^{\*}\tilde{\Psi}\_{t}^{-1}\tilde{\pmb{\eta}}\_{t}+\cdots
$$
\n(138)

Therefore, the backward message is approximately Gaussian with the mean and covariance given by

$$
\beta(\mathbf{x}_{t-1}) \approx \mathcal{N}(\mathbf{x}_{t-1}|\boldsymbol{\eta}_{t-1}, \Psi_{t-1}), \qquad (139)
$$

$$
\Psi_{t-1}^{-1} = \langle A^T A \rangle - \langle A \rangle^T \Psi_t^* \langle A \rangle, \tag{140}
$$

$$
\eta_{t-1} = \Psi_{t-1} < A >^T \Psi_t^* \tilde{\Psi}_t^{-1} \tilde{\eta}_t = \Psi_{t-1} < A >^T (I + \tilde{\Psi}_t)^{-1} \tilde{\eta}_t. \tag{141}
$$

#### Computing marginals of latent variables using  $\alpha$  and  $\beta$

Using the  $\alpha$  and  $\beta$  recursions in the forward/backward algorithm, we can compute the marginals of the latent variables.

$$
p(\mathbf{x}_t|\mathbf{y}_{1:T}) = p(\mathbf{x}_t|\mathbf{y}_{1:t}, \mathbf{y}_{t+1:T}),
$$
\n(142)

$$
\propto p(\mathbf{y}_{t+1:T}|\mathbf{x}_t, \mathbf{y}_{1:t})p(\mathbf{x}_t|\mathbf{y}_{1:t}) = p(\mathbf{y}_{t+1:T}|\mathbf{x}_t)p(\mathbf{x}_t|\mathbf{y}_{1:t}) = \beta(\mathbf{x}_t)\alpha(\mathbf{x}_t),
$$
\n(143)

$$
\propto \mathcal{N}(\mathbf{x}_t|\boldsymbol{\omega}_t, \Upsilon_t) \tag{144}
$$

where

$$
\Upsilon_t^{-1} = \Psi_t^{-1} + \Sigma_t^{-1}, \tag{145}
$$

$$
\omega_t = \Upsilon_t (\Psi_t^{-1} \eta_t + \Sigma_t^{-1} \mu_t). \tag{146}
$$

We also need to compute pairwise marginals of latent variables, given by

$$
p(\mathbf{x}_{t}, \mathbf{x}_{t+1} | \mathbf{y}_{1:T}) = p(\mathbf{x}_{t}, \mathbf{x}_{t+1} | \mathbf{y}_{1:t}, \mathbf{y}_{t+1}, \mathbf{y}_{t+2:T}),
$$
  
\n
$$
\propto p(\mathbf{y}_{t+1}, \mathbf{y}_{t+2:T} | \mathbf{x}_{t}, \mathbf{x}_{t+1}, \mathbf{y}_{1:t}) p(\mathbf{x}_{t+1} | \mathbf{x}_{t}, \mathbf{y}_{1:t}) p(\mathbf{x}_{t} | \mathbf{y}_{1:t}),
$$
  
\n
$$
\propto p(\mathbf{y}_{t+1} | \mathbf{x}_{t+1}) p(\mathbf{y}_{t+2:T} | \mathbf{x}_{t+1}) p(\mathbf{x}_{t+1} | \mathbf{x}_{t}) p(\mathbf{x}_{t} | \mathbf{y}_{1:t}),
$$
  
\n
$$
\propto \beta(\mathbf{x}_{t+1}) \exp \left( \langle \log(p(\mathbf{y}_{t+1} | \mathbf{x}_{t+1}) p(\mathbf{x}_{t+1} | \mathbf{x}_{t})) \rangle_{q_{\theta}(\theta)} \right) \alpha(\mathbf{x}_{t}),
$$
\n(147)

which are jointly Gaussian

$$
p\begin{pmatrix} \mathbf{x}_t \\ \mathbf{x}_{t+1} \end{pmatrix} = \mathcal{N}\left( \begin{bmatrix} \boldsymbol{\omega}_t \\ \boldsymbol{\omega}_{t+1} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Upsilon}_t & \boldsymbol{\Upsilon}_{t,t+1} \\ \boldsymbol{\Upsilon}_{t,t+1}^T & \boldsymbol{\Upsilon}_{t+1}^T \end{bmatrix} \right).
$$
(148)

,

To compute the cross-covariance  $\Upsilon_{t,t+1}$ , we first compute the second derivatives of log of eq. 147 w.r.t.  $[\mathbf{x}_t \ \mathbf{x}_{t+1}]^T$ :

$$
\frac{\partial^2 \log \int d\theta q_{\theta}(\theta) p(\mathbf{x}_t, \mathbf{x}_{t+1} | \mathbf{y}_{1:T})}{\partial [\mathbf{x}_t | \mathbf{x}_{t+1}]^2} = -\begin{bmatrix} \Sigma_t^{*-1} & -\langle A \rangle^{\top} \\ -\langle A \rangle & \Psi_{t+1}^{-1} + I + W_{t+1} \end{bmatrix},\tag{149}
$$

where

$$
W_{t+1} = \frac{\partial^2}{\partial \mathbf{x}_{t+1}^2} < \log p(\mathbf{y}_{t+1}|\mathbf{x}_{t+1}) >_{q(\theta)},\tag{150}
$$

$$
= \sum_{s=1}^p \left[ \Sigma_{\mathbf{c}_s} + (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_{t+1}) (\boldsymbol{\mu}_{\mathbf{c}_s} + \Sigma_{\mathbf{c}_s} \mathbf{x}_{t+1})^T \right] e^{\mathbf{x}_{t+1}^T \boldsymbol{\mu}_{\mathbf{c}_s} + \frac{1}{2} \mathbf{x}_{t+1}^T \Sigma_{\mathbf{c}_s} \mathbf{x}_{t+1} + \mathbf{d}_s}, \tag{151}
$$

evaluated at  $\mathbf{x}_{t+1} = \boldsymbol{\omega}_{t+1}$ . By negating and inverting the matrix in eq. 149, and using the *Schur* complement, we can obtain  $\Upsilon_{t,t+1}$ 

$$
\Upsilon_{t,t+1} = -(\Sigma_t^{*-1} - \langle A \rangle^T (\Psi_{t+1}^{-1} + I + W_{t+1})^{-1} \langle A \rangle)^{-1} (-\langle A \rangle^T) (\Psi_{t+1}^{-1} + I + W_{t+1})^{-1}.
$$
 (152)

### Computing sufficient statistics of latent variables

Using  $q_{\mathbf{x}}(\mathbf{x}_{0:1}^{(i)})$  $\binom{0,1}{0,T}$ , we can compute the sufficient statistics of latent variables (that are used in M step).

$$
W_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \mathbf{x}_{t-1}^{(i)} \rangle = \sum_{t=1}^{T} \Upsilon_{t-1}^{(i)} + \omega_{t-1}^{(i)} \omega_{t-1}^{(i)} , \qquad S_{A^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t-1}^{(i)} \mathbf{x}_{t}^{(i)} \rangle = \sum_{t=1}^{T} \Upsilon_{t-1,t}^{(i)} + \omega_{t-1}^{(i)} \omega_{t}^{(i)} , \tag{153}
$$
\n
$$
S_{C^{(i)}} = \sum_{t=1}^{T} \langle \mathbf{x}_{t}^{(i)} \rangle \mathbf{y}_{t}^{(i)} \rangle = \sum_{t=1}^{T} \omega_{t}^{(i)} \mathbf{y}_{t}^{(i)} \rangle.
$$

### Hyperaparameter estimation

We take the derivatives of the variational lower bound w.r.t. each hyperparameter to obtain update rules. The lower bound is simplified as below:

$$
\log p(\mathbf{y}_{1:T}^{(1:r)}) \geq \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log \frac{p(\theta, \mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)})}{q(\theta, \mathbf{x}_{0:T}^{(1:r)})},
$$
\n
$$
= \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log p(\mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)}) |\theta) - \int d\mathbf{x}_{0:T}^{(1:r)} q(\mathbf{x}_{0:T}^{(1:r)}) \log q(\mathbf{x}_{0:T}^{(1:r)}) + \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log \frac{p(\theta)}{q(\theta)},
$$
\n
$$
= \sum_{i=1}^{r} \log Z'_{(i)} + \int d\theta \, d\mathbf{x}_{0:T}^{(1:r)} q(\theta, \mathbf{x}_{0:T}^{(1:r)}) \log \frac{p(\theta)}{q(\theta)},
$$
\n(155)

where the line is true from eq. 114, i.e.,

$$
-\int d\mathbf{x}_{0:T}^{(1:r)} q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)}) \log q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)}) = -\int d\mathbf{x}_{0:T}^{(1:r)} q_{\mathbf{x}}(\mathbf{x}_{0:T}^{(1:r)}) \mathbb{E}_{q_{\theta}(\theta)} \log p(\mathbf{x}_{0:T}^{(1:r)}, \mathbf{y}_{1:T}^{(1:r)} | \theta) + \sum_{i=1}^{r} \log Z'_{(i)}.
$$
 (156)

So, we need to consider the second term in RHS of eq. 155 for hyperparameter update (the integration w.r.t. x is omitted, since the integrand is independent of  $x$ )

$$
\int d\mathbf{a}^{(1:r)} dC q(\mathbf{a}^{(1:r)}) q(C) \log \frac{p(\mathbf{a}^{(1:r)}, C)}{q(\mathbf{a}^{(1:r)}) q(C)} = -KL(C) - KL(\mathbf{a}^{(1:r)}),
$$
\n(157)

The first term,  $KL(C)$  is given by <sup>1</sup>

$$
KL(C) = \int dC q_C(C) \log \frac{q_C(C)}{p(C|\gamma)},
$$
  
\n
$$
= \sum_{s=1}^p \int d\mathbf{c}_s \, \mathcal{N}(\mathbf{c}_s | \boldsymbol{\mu}_{\mathbf{c}_s}, \Sigma_{\mathbf{c}_s}) \log \frac{\mathcal{N}(\mathbf{c}_s | \boldsymbol{\mu}_{\mathbf{c}_s}, \Sigma_{\mathbf{c}_s})}{\mathcal{N}(\mathbf{c}_s | \mathbf{0}, \gamma^{-1} I)},
$$
  
\n
$$
= \sum_{s=1}^p \left(-\frac{1}{2} \log |\gamma \Sigma_{\mathbf{c}_s}| + \frac{1}{2} \text{Tr} \left[\gamma (\Sigma_{\mathbf{c}_s} - \gamma^{-1} I + \boldsymbol{\mu}_{\mathbf{c}_s} \boldsymbol{\mu}_{\mathbf{c}_s}^T) \right] \right).
$$
(159)

The first derivative expression w.r.t.  $\gamma$  gives us the following update:

$$
\gamma^{-1} = \frac{1}{p} \sum_{s=1}^{p} \text{Tr}[\Sigma_{\mathbf{c}_s} + \boldsymbol{\mu}_{\mathbf{c}_s} \boldsymbol{\mu}_{\mathbf{c}_s}^T], \qquad (160)
$$

Similarly, the second term is given by

$$
KL(\mathbf{a}^{(1:r)}) = \int d\mathbf{a}^{(1:r)} q_{\mathbf{a}}(\mathbf{a}^{(1:r)}) \log \frac{q_{\mathbf{a}}(\mathbf{a}^{(1:r)})}{p(\mathbf{a}^{(1:r)}|\bar{\mathbf{a}}, \sigma^2, \tau^2)},
$$
\n(161)

$$
= \int d\mathbf{a}^{(1:r)} \mathcal{N}(\mathbf{a}^{(1:r)} | \boldsymbol{\mu}_{\mathbf{a}}, \Sigma_{\mathbf{a}}) \log \frac{\mathcal{N}(\mathbf{a} | \boldsymbol{\mu}_{\mathbf{a}}, \Sigma_{\mathbf{a}})}{\mathcal{N}(\mathbf{a}^{(1:r)} | \mathbf{\bar{a}}^{(1:r)}, K)},
$$
\n(162)

$$
= -\frac{1}{2}\log|K^{-1}\Sigma_{\mathbf{a}}| + \frac{1}{2}\text{Tr}\left[K^{-1}\Sigma_{\mathbf{a}}\right] + \frac{1}{2}(\boldsymbol{\mu}_{\mathbf{a}} - \bar{\mathbf{a}}^{(1:r)})^{\top}K^{-1}(\boldsymbol{\mu}_{\mathbf{a}} - \bar{\mathbf{a}}^{(1:r)}) + const.
$$
(163)

The first derivative w.r.t.  $\bar{a}$  is given by

$$
\frac{\partial}{\partial \bar{\mathbf{a}}} KL(\mathbf{a}^{(1:r)}) = \frac{1}{2} \frac{\partial}{\partial \bar{\mathbf{a}}} (\boldsymbol{\mu}_\mathbf{a} - \bar{\mathbf{a}}^{(1:r)})^\top K^{-1} (\boldsymbol{\mu}_\mathbf{a} - \bar{\mathbf{a}}^{(1:r)}), \tag{164}
$$

$$
= \frac{1}{2} \frac{\partial}{\partial \bar{\mathbf{a}}} (\boldsymbol{\mu}_\mathbf{a} - E \bar{\mathbf{a}})^{\top} K^{-1} (\boldsymbol{\mu}_\mathbf{a} - E \bar{\mathbf{a}})
$$
(165)

where  $E = \mathbf{1}_r \otimes I_{k^2}$ , and this gives us the update rule:

$$
\bar{\mathbf{a}} = (E^{\top} K^{-1} E)^{-1} (E^{\top} K^{-1} \mu_{\mathbf{a}}). \tag{166}
$$

The first derivative w.r.t. kernel parameters (denoted by  $\alpha = \{\sigma^2, \tau^2\}$ ) is given by

$$
\frac{\partial}{\partial \alpha} KL(\mathbf{a}^{(1:r)}) = \frac{1}{2} \text{Tr} \left( K^{-1} \frac{\partial K}{\partial \alpha} \right) - \frac{1}{2} \text{Tr} \left( K^{-1} \frac{\partial K}{\partial \alpha} K^{-1} (\Sigma_{\mathbf{a}} + (\mu_{\mathbf{a}} - \bar{\mathbf{a}}^{(1:r)}) (\mu_{\mathbf{a}} - \bar{\mathbf{a}}^{(1:r)})^\top) \right), \tag{167}
$$

$$
= \frac{1}{2}\text{Tr}\left(K^{-1}\frac{\partial K}{\partial \alpha}(I - K^{-1}(\Sigma_{\mathbf{a}} + (\boldsymbol{\mu}_{\mathbf{a}} - \bar{\mathbf{a}}^{(1:r)})(\boldsymbol{\mu}_{\mathbf{a}} - \bar{\mathbf{a}}^{(1:r)})^{\top}))\right),\tag{168}
$$

where the first derivative of  $K(i, j)$  w.r.t.  $\alpha$  is given by

$$
\frac{\partial}{\partial \tau^2} K(i,j) = \frac{1}{2\tau^4} (i-j)^2 (\sigma^2 + \epsilon \delta_{ij}) \exp\left(-\frac{1}{2\tau^2} (i-j)^2\right) I_{k^2} = \frac{1}{2\tau^4} (i-j)^2 K(i,j),\tag{169}
$$

$$
\frac{\partial}{\partial \sigma^2} K(i,j) = \exp\left(-\frac{1}{2\tau^2}(i-j)^2\right) I_{k^2}.
$$
\n(170)

We update  $\alpha$  numerically using the derivative expression above.

$$
KL(\tilde{\boldsymbol{\mu}}, \tilde{\Sigma}||\boldsymbol{\mu}, \Sigma) = -\frac{1}{2}\log|\tilde{\Sigma}\Sigma^{-1}| + \frac{1}{2}\text{Tr}\left[\Sigma^{-1}(\tilde{\Sigma} - \Sigma + (\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu})(\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu})^T)\right].
$$
\n(158)

 $1$ <sup>1</sup>The formula of KL divergence between two Gaussians is given by:

# Illustration with simulated data



Figure 2. Illustration of non-stationarity in population dynamics (data simulated from Model II). A: Raster plots of spontaneous activity from 40 neurons during 10 seconds of recording for simulated trials 1 and 100. We assumed that the two sub-populations (blue and red) have negative correlation at trial 1 and positive correlation at trial 100. B: Recovered correlations. Our Model II (red) accurately recovers the correlations between two groups across trials (RMSE: 0.04), while other methods perform poorly: independent PLDSs fit to each trial individually give noisy results (RMSE 0.06) and a single PLDS fit across all trials cannot capture the change in correlation (RMSE 0.44). C: Estimation of off-diagonal in dynamics matrices. We fixed the loading matrix  $C$  to its true value to avoid issues with non-identifiability of parameters in LDS models. The off-diagonal term  $A_{12}$  estimated by our model matched the true values well, whereas the independent PLDS produced noisy estimates, and the fixed PLDS cannot capture the change in  $A_{12}$ .

We tested Model II using a simulation of spontaneous activity from a population of 40 neurons (simulated from Model II). We assumed that the population could be split into two sub-populations of size 20 neurons each, and simulated an experiment in which the correlation across the two sub-populations changed dramatically across the experiment: Specifically, we generated a 2-d latent state that controls correlations in firing rates between the two groups of neurons, and adjusted the off-diagonal term in the dynamics matrix  $(A_{12})$  such that the correlation between the groups varied slowly from  $-1$  to 1 across 100 trials, where the length of each trial is  $T = 200$ . Other elements of A were adjusted such that the stationary covariance of the system was kept constant.

We fit Model II N-PLDS, a single PLDS, and 100 independent PLDSs to the data. Our model accurately recovered the correlation change in z across trials, while the single PLDS was not able to capture the non-stationarity and the independent PLDSs exhibited noisy correlations (Fig. 2). Finally, our model also accurately recovered the off-diagonal parameter  $A_{12}$ (Fig. 2 C). For panel C only, we set the loading matrix C to the ground truth value for each of the models (Model II, fixed PLDS, separate PLDSs). LDS models suffer from non-identifiability of parameters, implying that estimated parameters do not necessarily match the true parameters even for perfect model fits.

# Illustration with real data

Finally, we analyzed a dataset of spontaneous activity recorded from a population of 40 neurons from macaque visual cortex. The details of data collection are described in [1] and the data is available from [2]. Using the spike-sorting information provided in the dataset, we selected the spike-cluster with highest signal-to-noise ratio from each recording channel, and out of those 46 units kept the 40 units with highest firing rates. As the original data consisted of one continuous recording of length 15 minutes, we divided the data into 30 'epochs' of length 30 seconds each, and used every 5th epoch (20% of the data) for testing and the rest (80% of data) for training.

In this data, the mean firing rates are almost constant across time, while the correlations increase at the end of the experiment (Fig. 3 A). After estimating the parameters of our N-PLDS (Model II) from the training data, we computed the predictive distribution on the dynamics matrices A<sup>∗</sup> for the test data. Using these parameters, we drew samples for spikes to compute the mean firing rates for each trial (Fig. 3 A), as well as the mean pairwise cross-correlations across all neuron pairs. The correlations estimated from N-PLDS (Model II) matched those in the data. For PLDS with fixed parameters, the estimated firing rates and correlations are constant across epochs (Fig. 3 B). To quantify these results, we computed the RMSE in the prediction of mean firing rates and mean correlations on test epochs. The RMSEs on mean firing rate estimation for PLDS are 0.0156, 0.0182, 0.0188 for  $k = 1, 2, 4$ , respectively, while RMSE of N-PLDS is 0.0080 ( $k = 4$ ). The RMSE on mean correlation estimation in PLDSs is 0.0138 (same for  $k = 1, 2, 4$ ) and 0.0087 ( $k = 4$ ) in N-PLDS.



Figure 3. Non-stationary population dynamics (data from [1]). A: Summary statistics of samples from N-PLDS (Model II) with non-stationarity dynamics matrix A for different dimensions of latent dynamics  $(k = 1, 2, 4)$ . The top plot shows the mean firing rate of 40 neurons during 30 epochs, showing that there is only a slight systematic drift in mean firing rate. Each dot represents predicted mean firing rates for the held-out data (6 trials). The bottom plot shows the mean correlation of the spike counts. All three N-PLDS models capture the increase in correlation at the end of the experiment, with the  $k = 4$  capturing it most accurately. B: Comparison to using a PLDS model with fixed parameters  $(k = 1, 2, 4)$ . Both the mean firing rate and correlation in PLDS are constant across epochs. As a consequence, the best RMSE on mean correlation estimation in PLDS is 0.0138  $(k = 1)$  compared to 0.0087  $(k = 4)$  in N-PLDS.

## References

- 1. Cheng C.J. Chu, Ping F. Chien, and Chou P. Hung. Tuning dissimilarity explains short distance decline of spontaneous spike correlation in macaque  $\{V1\}$ . Vision Research, 96(0):113 – 132, 2014.
- 2. Chou P. Hung Cheng C. J. Chu, Ping F. Chien. Multi-electrode recordings of ongoing activity and responses to parametric stimuli in macaque v1. CRCNS.org., 2014.