

Appendix S1

In the following, we list the gradients necessary for implementing our model. For any parameter θ , the derivative of the conditional log-likelihood (Equation 3, main text) is

$$\frac{\partial}{\partial \theta} \log p(y | x, \theta) = \sum_{c,s} p(c, s | x, y) \frac{\partial}{\partial \theta} \log p(c, s, y | x) \quad (1)$$

where

$$\frac{\partial}{\partial \theta} \log p(c, s, y | x) = \frac{\partial}{\partial \theta} \log p(c, s | x) + \frac{\partial}{\partial \theta} \log p(y | c, s, x). \quad (2)$$

The first term can be written as

$$\frac{\partial}{\partial \theta} \log p(c, s | x) = \frac{\partial}{\partial \theta} \log \frac{f(c, s; x)}{\sum_{c,s} f(c, s; x)} \quad (3)$$

$$= \frac{\partial}{\partial \theta} \log p(c, s, x) - \sum_{c,s} p(c, s | x) \frac{\partial}{\partial \theta} \log f(c, s; x), \quad (4)$$

where $f(c, s; x)$ (the *gate*) is defined as in Equation 4 of the main text to be the unnormalized conditional probability

$$f(c, s; x) = |\lambda_{cs}|^{\frac{N}{2}} |L_c| \exp\left(-\frac{1}{2} x^\top \lambda_{cs} L_c L_c^\top x\right). \quad (5)$$

We replaced the precision matrices K_c and M_c (Equations 4 and 5, main text) by their Cholesky decompositions $L_c L_c^\top$ and $B_c B_c^\top$, respectively, and optimized over the Cholesky factors instead. The conditional density $p(y | c, s, x)$ (the *expert*) is Gaussian and therefore easily normalized,

$$p(y | x, c, s) = \frac{\lambda_{cs}^{\frac{M}{2}}}{(2\pi)^{\frac{M}{2}}} |B_c| \exp\left(-\frac{1}{2} (y - A_c x)^\top \lambda_{cs} B_c B_c^\top (y - A_c x)\right). \quad (6)$$

What is left is to compute the gradients of $\log p(y | x, c, s)$ and $\log f(c, s; x)$,

$$\frac{\partial}{\partial L_c} \log f(c, s; x) = \text{diag}(L_c)^{-1} - \lambda_{cs} x x^\top L_c, \quad (7)$$

$$\frac{\partial}{\partial \lambda_{cs}} \log f(c, s; x) = \frac{N}{2} \lambda_{cs}^{-1} - \frac{1}{2} x^\top L_c L_c^\top x, \quad (8)$$

$$\frac{\partial}{\partial \lambda_{cs}} \log p(y | c, s, x) = \frac{M}{2} \lambda_{cs}^{-1} - \frac{1}{2} (y - A_c x)^\top B_c B_c^\top (y - A_c x), \quad (9)$$

$$\frac{\partial}{\partial B_c} \log p(y | c, s, x) = \text{diag}(B_c)^{-1} - \lambda_{cs} (y - A_c x)(y - A_c x)^\top B_c, \quad (10)$$

$$\frac{\partial}{\partial A_c} \log p(y | c, s, x) = \lambda_{cs} B_c B_c^\top (y - A_c x) x^\top. \quad (11)$$

Here, M is the dimensionality of y ($M = 3$ for the multiscale MCGSM and $M = 1$ for the regular MCGSM) and N is the dimensionality of the neighborhood x . The function diag returns a matrix with all off-diagonals set to zero, $\text{diag}(M)_{ij} = \delta_{ij} M_{ij}$.

We tried regularizing the parameters by specifying different prior distributions for the parameters and optimizing the posterior probability of the parameters instead of their likelihood. However, we found no improvement in the model's performance. We also tried adding means to the gates and experts, which also did not help to improve the performance of the model when evaluated on natural images.