

Appendix A: derivation for the marginal distribution in eq. 6

First, we show a general derivation for the marginal distribution over \mathbf{F} . First, we concatenate $\{\tilde{\mathbf{y}}_t\}_{t=1}^T$ vertically to get a big vector $\hat{\mathbf{y}} \in \mathbb{R}^{(TDN) \times 1}$ resulting in the following distribution

$$\hat{\mathbf{y}}|\tilde{\mathbf{f}} \sim \mathcal{N}(\mathbf{h} \otimes \tilde{\mathbf{f}}, \mathbf{I}_T \otimes \Delta) \quad (16)$$

where \mathbf{h} is a $T \times 1$ vector of ones. We can replace \mathbf{h} with a $(TDN) \times (DN)$ matrix \mathbf{H} such that $\mathbf{H}\tilde{\mathbf{f}} = \mathbf{h} \otimes \tilde{\mathbf{f}}$. The marginal distribution of $\hat{\mathbf{y}}$ can be written as

$$\begin{aligned} p(\hat{\mathbf{y}}|\mathbf{K}) &= \int \mathcal{N}(\hat{\mathbf{y}}|\mathbf{H}\tilde{\mathbf{f}}, \mathbf{I}_T \otimes \Delta) \mathcal{N}(\tilde{\mathbf{f}}|\mathbf{0}, \mathbf{I}_N \otimes \mathbf{K}) d\tilde{\mathbf{f}} \\ &= \frac{1}{Z} \int \exp \left(-\frac{1}{2}(\hat{\mathbf{y}} - \mathbf{H}\tilde{\mathbf{f}})^\top (\mathbf{I}_T \otimes \Delta)^{-1} (\hat{\mathbf{y}} - \mathbf{H}\tilde{\mathbf{f}}) - \frac{1}{2}\tilde{\mathbf{f}}^\top (\mathbf{I}_N \otimes \mathbf{K})^{-1} \tilde{\mathbf{f}} \right) d\tilde{\mathbf{f}} \\ &= \frac{1}{Z'} \exp \left(-\frac{1}{2}\hat{\mathbf{y}}^\top (\mathbf{I}_T \otimes \Delta + \mathbf{H}(\mathbf{I}_N \otimes \mathbf{K})\mathbf{H}^\top)^{-1} \hat{\mathbf{y}} \right) \\ &= \mathcal{N}(\hat{\mathbf{y}}|\mathbf{0}, \mathbf{I}_T \otimes \Delta + \mathbf{H}(\mathbf{I}_N \otimes \mathbf{K})\mathbf{H}^\top) \end{aligned} \quad (17)$$

The new covariance $\mathbf{I}_T \otimes \Delta + \mathbf{H}(\mathbf{I}_N \otimes \mathbf{K})\mathbf{H}^\top$ is a $(TDN) \times (TDN)$ matrix which is computationally not invertible in practice. However the heavy inversion can be resolved by applying matrix inversion lemma and the property of Kronecker product when calculating the log-likelihood.

Here, we provide another way of marginalizing out \mathbf{F} which consists of multiple multivariate normal distributions with smaller scale covariance matrices. Instead of vectorizing $\tilde{\mathbf{Y}}$ matrix into $\hat{\mathbf{y}}$ and dealing with one multivariate normal with a big covariance matrix, we work on the integration with the Gaussian distribution for data in eq. 5. The marginal distribution can be written as

$$p(\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_T|\mathbf{K}) = \int \mathcal{N}(\tilde{\mathbf{f}}|\mathbf{0}, \mathbf{I}_N \otimes \mathbf{K}) \prod_{t=1}^T \mathcal{N}(\tilde{\mathbf{y}}_t|\tilde{\mathbf{f}}, \Delta) d\tilde{\mathbf{f}} \quad (18)$$

Given a set of data observations $\{\tilde{\mathbf{y}}_t\}_{t=1}^T$, we can write the probability density function of $\mathcal{N}(\tilde{\mathbf{y}}_t|\tilde{\mathbf{f}}, \Delta)$ as $\mathcal{N}(\tilde{\mathbf{f}}|\tilde{\mathbf{y}}_t, \Delta)$ which is just an exponential function of a negative quadratic function. According to the property of the product of Gaussian densities, let $\mathcal{N}_x(m, \Sigma)$ denote a density of x , then

$$\begin{aligned} \mathcal{N}_x(m_1, \Sigma_1) \mathcal{N}_x(m_2, \Sigma_2) &= c_c \mathcal{N}_x(m_c, \Sigma_c), & c_c &= \mathcal{N}_{m_1}(m_2, \Sigma_1 + \Sigma_2), \\ m_c &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1} (\Sigma_1^{-1} m_1 + \Sigma_2^{-1} m_2), & \Sigma_c &= (\Sigma_1^{-1} + \Sigma_2^{-1})^{-1}. \end{aligned} \quad (19)$$

We can apply the property to the integration in eq. 18 in a chain style from $\mathcal{N}(\tilde{\mathbf{f}}|\tilde{\mathbf{y}}_1, \Delta)$ all the way to $\mathcal{N}(\tilde{\mathbf{f}}|\mathbf{0}, \mathbf{I}_N \otimes \mathbf{K})$:

$$\begin{aligned}
(1). \quad & P_1 = \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_1, \Delta) \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_2, \Delta) = c_1 \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{2}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2), \frac{1}{2}\Delta \right), \\
& c_1 = \mathcal{N}_{\tilde{\mathbf{y}}_1}(\tilde{\mathbf{y}}_2, 2\Delta), \\
(2). \quad & P_2 = c_1 \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{2}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2), \frac{1}{2}\Delta \right) \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_3, \Delta) = c_1 c_2 \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{3}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2 + \tilde{\mathbf{y}}_3), \frac{1}{3}\Delta \right), \\
& c_2 = \mathcal{N}_{\frac{1}{2}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2)}(\tilde{\mathbf{y}}_3, \frac{3}{2}\Delta), \\
(3). \quad & P_3 = c_1 c_2 \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{3}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2 + \tilde{\mathbf{y}}_3), \frac{1}{3}\Delta \right) \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_4, \Delta) = c_1 c_2 c_3 \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{4}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2 + \tilde{\mathbf{y}}_3 + \tilde{\mathbf{y}}_4), \frac{1}{4}\Delta \right), \\
& c_3 = \mathcal{N}_{\frac{1}{3}(\tilde{\mathbf{y}}_1 + \tilde{\mathbf{y}}_2 + \tilde{\mathbf{y}}_3)}(\tilde{\mathbf{y}}_4, \frac{4}{3}\Delta), \\
& \vdots \\
(T-1). \quad & P_{T-1} = \prod_{t=1}^{T-2} c_t \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{T-1} \sum_{t=1}^{T-1} \tilde{\mathbf{y}}_t, \frac{1}{T-1}\Delta \right) \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_T, \Delta) \prod_{t=1}^{T-1} c_t \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{T} \sum_{t=1}^T \tilde{\mathbf{y}}_t, \frac{1}{T}\Delta \right), \\
& c_{T-1} = \mathcal{N}_{\frac{1}{T-1} \sum_{t=1}^{T-1} \tilde{\mathbf{y}}_t}(\tilde{\mathbf{y}}_T, \frac{T}{T-1}\Delta), \\
(T). \quad & P_T = \prod_{t=1}^{T-1} c_t \mathcal{N}_{\tilde{\mathbf{f}}} \left(\frac{1}{T} \sum_{t=1}^T \tilde{\mathbf{y}}_t, \frac{1}{T}\Delta \right) \mathcal{N}_{\tilde{\mathbf{f}}}(\mathbf{0}, \mathbf{I}_N \otimes \mathbf{K}) = \prod_{t=1}^T c_t \mathcal{N}_{\tilde{\mathbf{f}}}(\cdot, \cdot), \\
& c_T = \mathcal{N}_{\frac{1}{T} \sum_{t=1}^T \tilde{\mathbf{y}}_t}(\mathbf{0}, \frac{1}{T}\Delta + \mathbf{I}_N \otimes \mathbf{K})
\end{aligned}$$

Therefore, we can write eq. [18](#) as

$$\begin{aligned}
p(\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_T | \mathbf{K}) &= \int \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_1, \Delta) \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_2, \Delta) \dots \mathcal{N}_{\tilde{\mathbf{f}}}(\tilde{\mathbf{y}}_T, \Delta) \mathcal{N}_{\tilde{\mathbf{f}}}(\mathbf{0}, \mathbf{I}_N \otimes \mathbf{K}) d\tilde{\mathbf{f}} = \int P_T d\tilde{\mathbf{f}} \\
&= \int \prod_{t=1}^T c_t \mathcal{N}_{\tilde{\mathbf{f}}}(\cdot, \cdot) d\tilde{\mathbf{f}} = \prod_{t=1}^T c_t \int \mathcal{N}_{\tilde{\mathbf{f}}}(\cdot, \cdot) d\tilde{\mathbf{f}} = \prod_{t=1}^T c_t
\end{aligned}$$

Its log likelihood is

$$\log p(\tilde{\mathbf{y}}_1, \tilde{\mathbf{y}}_2, \dots, \tilde{\mathbf{y}}_T) = \sum_{t=1}^T \log c_t \quad (20)$$

$$= \sum_{t=1}^{T-1} \left[-\frac{1}{2} \log \left| \frac{t+1}{t} \Delta \right| - \frac{1}{2} \left(\tilde{\mathbf{y}}_{t+1} - \frac{1}{t} \sum_{j=1}^t \tilde{\mathbf{y}}_j \right)^\top \left(\frac{t+1}{t} \Delta \right)^{-1} \left(\tilde{\mathbf{y}}_{t+1} - \frac{1}{t} \sum_{j=1}^t \tilde{\mathbf{y}}_j \right) \right] \\ - \frac{1}{2} \log \left| \frac{1}{T} \Delta + \mathbf{I} \otimes \mathbf{K} \right| - \frac{1}{2T} \sum_{t=1}^T \tilde{\mathbf{y}}_t^\top \left(\frac{1}{T} \Delta + \mathbf{I} \otimes \mathbf{K} \right)^{-1} \frac{1}{T} \sum_{t=1}^T \tilde{\mathbf{y}}_t \quad (21)$$

$$= \sum_{t=1}^{T-1} \left[-\frac{DN}{2} \log \left(\frac{t+1}{t} \right) - \frac{1}{2} \log |\Delta| - \frac{1}{2} \left(\sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1} - \frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j \right)^\top \Delta^{-1} \left(\sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1} - \frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j \right) \right] \\ - \frac{1}{2} \log \left| \frac{1}{T} \Delta + \mathbf{I} \otimes \mathbf{K} \right| - \frac{1}{2T} \sum_{t=1}^T \tilde{\mathbf{y}}_t^\top \left(\frac{1}{T} \Delta + \mathbf{I} \otimes \mathbf{K} \right)^{-1} \frac{1}{T} \sum_{t=1}^T \tilde{\mathbf{y}}_t \quad (22)$$

$$= -\frac{DN}{2} \log(T) + \sum_{t=1}^{T-1} \left[-\frac{1}{2} \log |\Delta| - \frac{1}{2} \left(\sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1} - \frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j \right)^\top \Delta^{-1} \left(\sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1} - \frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j \right) \right] \\ - \frac{1}{2} \log \left| \frac{1}{T} \Delta + \mathbf{I} \otimes \mathbf{K} \right| - \frac{1}{2T} \sum_{t=1}^T \tilde{\mathbf{y}}_t^\top \left(\frac{1}{T} \Delta + \mathbf{I} \otimes \mathbf{K} \right)^{-1} \frac{1}{T} \sum_{t=1}^T \tilde{\mathbf{y}}_t \quad (23)$$

$$= \sum_{t=1}^{T-1} \log \mathcal{N} \left(\sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1} \middle| \frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j, \Delta \right) - \frac{1}{2\sqrt{T}} \sum_{t=1}^T \tilde{\mathbf{y}}_t^\top (\Delta + T\mathbf{I} \otimes \mathbf{K})^{-1} \frac{1}{\sqrt{T}} \sum_{t=1}^T \tilde{\mathbf{y}}_t - \frac{1}{2} \log |\Delta + T\mathbf{I} \otimes \mathbf{K}| \\ = \sum_{t=1}^{T-1} \log \mathcal{N} \left(\frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j \middle| \sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1}, \Delta \right) + \log \mathcal{N} \left(\frac{1}{\sqrt{T}} \sum_{j=1}^T \tilde{\mathbf{y}}_j \middle| \mathbf{0}, \Delta + T\mathbf{I} \otimes \mathbf{K} \right)$$

Thus the marginal distribution is

$$p(\tilde{\mathbf{y}}_1, \dots, \tilde{\mathbf{y}}_T | \mathbf{K}) = \mathcal{N} \left(\frac{1}{\sqrt{T}} \sum_{j=1}^T \tilde{\mathbf{y}}_j \middle| \mathbf{0}, \Delta + T\mathbf{I} \otimes \mathbf{K} \right) \prod_{t=1}^{T-1} \mathcal{N} \left(\frac{1}{\sqrt{t(t+1)}} \sum_{j=1}^t \tilde{\mathbf{y}}_j - \sqrt{\frac{t}{t+1}} \tilde{\mathbf{y}}_{t+1} \middle| \mathbf{0}, \Delta \right)$$

Appendix B: black box variational inference

The log marginal likelihood for eq. [7](#) can be lower bounded by introducing any distribution over latent variable which has support where true posterior $p(\mathbf{X} | \mathbf{Y}, \Delta, \theta)$ does, and then appealing to Jensen's inequality (due to the concavity of the logarithm function):

$$\log p(\mathbf{Y} | \Delta, \theta) = \log \int p(\mathbf{Y} | \mathbf{X}, \Delta, \theta) p(\mathbf{X}) d\mathbf{X} \geq \int q(\mathbf{X} | \lambda) \log \frac{p(\mathbf{Y} | \mathbf{X}, \Delta, \theta) p(\mathbf{X})}{q(\mathbf{X} | \lambda)} d\mathbf{X} \quad (24)$$

where $q(\mathbf{X} | \lambda)$ is the variational approximating distribution for the true posterior controlled by some free variational parameters λ . We assume $q(\mathbf{X} | \lambda)$ to be a standard normal distribution. In E-step, we optimize the Evidence Lower BOund (ELBO),

$$\mathcal{L}(\lambda) \triangleq \mathbb{E}_{q(\mathbf{X} | \lambda)} [\log p(\mathbf{Y} | \mathbf{X}, \Delta, \theta) + \log p(\mathbf{X}) - \log q(\mathbf{X} | \lambda)] \quad (25)$$

A standard gradient descent method can be used to maximize the ELBO over the variational parameter with analytic computation of the expectation. However, the expectation of the first term in eq. [25](#) doesn't have a closed-form solution. Therefore, we will employ the Black Box Variational Inference (BBVI) [\[20\]](#) to maximize the ELBO with stochastic optimization. The BBVI algorithm requires the computation of noisy unbiased gradients of the ELBO with Monte Carlo samples from the variational

distribution,

$$\nabla_{\lambda} \mathcal{L}(\lambda) \approx \frac{1}{l} \sum_{i=1}^l \nabla_{\lambda} \log q(\mathbf{X}_l | \lambda) (\log p(\mathbf{Y} | \mathbf{X}_l, \Delta, \theta) + \log p(\mathbf{X}_l) - \log q(\mathbf{X}_l | \lambda)), \text{ where } \mathbf{X}_l \sim q(\mathbf{X} | \lambda). \quad (26)$$

This gradient involves calculating the log likelihood of $p(\mathbf{Y} | \mathbf{X}, \Delta, \theta)$ with $(DN) \times (DN)$ covariance matrices, which is the computational bottleneck of the evaluation. However, we can efficiently evaluate it with the nice property of Kronecker product.

Appendix C: inverting the covariance matrix

The key problem is to invert the covariance matrix $\mathbf{C} = T\mathbf{I}_N \otimes \mathbf{K} + \Sigma_N \otimes \Sigma_D$.

Let $\Sigma_D = \mathbf{U}_D \Lambda_D \mathbf{U}_D^{\top}$ and $\Sigma_N = \mathbf{U}_N \Lambda_N \mathbf{U}_N^{\top}$ be the eigen-decompositions of Σ_D and Σ_N . The covariance matrix \mathbf{C} can be factorized as

$$\begin{aligned} \mathbf{C} &= T\mathbf{I}_N \otimes \mathbf{K} + \Sigma_N \otimes \Sigma_D \\ &= T\mathbf{I}_N \otimes \mathbf{K} + (\mathbf{U}_N \Lambda_N \mathbf{U}_N^{\top}) \otimes (\mathbf{U}_D \Lambda_D \mathbf{U}_D^{\top}) \\ &= T\mathbf{I}_N \otimes \mathbf{K} + (\mathbf{U}_N \Lambda_N^{\frac{1}{2}} \Lambda_N^{\frac{1}{2}} \mathbf{U}_N^{\top}) \otimes (\mathbf{U}_D \Lambda_D^{\frac{1}{2}} \Lambda_D^{\frac{1}{2}} \mathbf{U}_D^{\top}) \\ &= T\mathbf{I}_N \otimes \mathbf{K} + \left(\mathbf{U}_N \Lambda_N^{\frac{1}{2}} \otimes \mathbf{U}_D \Lambda_D^{\frac{1}{2}} \right) \left(\Lambda_N^{\frac{1}{2}} \mathbf{U}_N^{\top} \otimes \Lambda_D^{\frac{1}{2}} \mathbf{U}_D^{\top} \right) \\ &= \left(\mathbf{U}_N \Lambda_N^{\frac{1}{2}} \otimes \mathbf{U}_D \Lambda_D^{\frac{1}{2}} \right) \left(\left(\mathbf{U}_N \Lambda_N^{\frac{1}{2}} \otimes \mathbf{U}_D \Lambda_D^{\frac{1}{2}} \right)^{-1} (T\mathbf{I}_N \otimes \mathbf{K}) \left(\Lambda_N^{\frac{1}{2}} \mathbf{U}_N^{\top} \otimes \Lambda_D^{\frac{1}{2}} \mathbf{U}_D^{\top} \right)^{-1} + \mathbf{I}_N \otimes \mathbf{I}_D \right) \\ &\quad \left(\Lambda_N^{\frac{1}{2}} \mathbf{U}_N^{\top} \otimes \Lambda_D^{\frac{1}{2}} \mathbf{U}_D^{\top} \right) \\ &= \left(\mathbf{U}_N \Lambda_N^{\frac{1}{2}} \otimes \mathbf{U}_D \Lambda_D^{\frac{1}{2}} \right) \left((T\Lambda_N^{-1}) \otimes (\Lambda_D^{-\frac{1}{2}} \mathbf{U}_D^{\top} \mathbf{K} \mathbf{U}_D \Lambda_D^{-\frac{1}{2}}) + \mathbf{I}_N \otimes \mathbf{I}_D \right) \left(\Lambda_N^{\frac{1}{2}} \mathbf{U}_N^{\top} \otimes \Lambda_D^{\frac{1}{2}} \mathbf{U}_D^{\top} \right). \quad (27) \end{aligned}$$

The complexity of inverting the first and the third terms in eq. 27 is $O(D^3 + N^3)$. The bottleneck is now inverting the second term in eq. 27. We define new notations $\tilde{\mathbf{K}} = \Lambda_D^{-\frac{1}{2}} \mathbf{U}_D^{\top} \mathbf{K} \mathbf{U}_D \Lambda_D^{-\frac{1}{2}}$ and $\tilde{\mathbf{C}} = T\Lambda_N^{-1} \otimes \tilde{\mathbf{K}} + \mathbf{I}_N \otimes \mathbf{I}_D$.

The problem is reduced to inverting the matrix $\tilde{\mathbf{C}}$. Therefore the second step is to exploit the compatibility of a Kronecker product plus a constant diagonal term with eigenvalue decomposition. Let $T\Lambda_N^{-1} = \mathbf{U}_T \Lambda_T \mathbf{U}_T^{\top}$ and $\tilde{\mathbf{K}} = \mathbf{U}_K \Lambda_K \mathbf{U}_K^{\top}$ be the eigen-decompositions of $T\Lambda_N^{-1}$ and $\tilde{\mathbf{K}}$. Thus,

$$\tilde{\mathbf{C}} = T\Lambda_N^{-1} \otimes \tilde{\mathbf{K}} + \mathbf{I}_N \otimes \mathbf{I}_D = (\mathbf{U}_T \otimes \mathbf{U}_K) (\Lambda_T \otimes \Lambda_K + \mathbf{I}_N \otimes \mathbf{I}_D) (\mathbf{U}_T^{\top} \otimes \mathbf{U}_K^{\top}), \quad (28)$$

Finally, combining eq. 27 and eq. 28 together to get

$$\mathbf{C} = \left(\mathbf{U}_N \Lambda_N^{\frac{1}{2}} \otimes \mathbf{U}_D \Lambda_D^{\frac{1}{2}} \right) (\mathbf{U}_T \otimes \mathbf{U}_K) (\Lambda_T \otimes \Lambda_K + \mathbf{I}_N \otimes \mathbf{I}_D) (\mathbf{U}_T^{\top} \otimes \mathbf{U}_K^{\top}) \left(\Lambda_N^{\frac{1}{2}} \mathbf{U}_N^{\top} \otimes \Lambda_D^{\frac{1}{2}} \mathbf{U}_D^{\top} \right). \quad (29)$$

Appendix D: more 2D latent representations for 22 odors

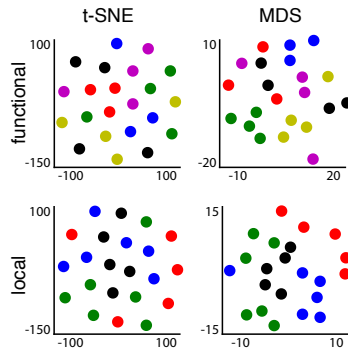


Figure 5: We analyzed the same dataset with t-SNE and MDS, and present the results obtained in the figures. Note that neither method is able to identify the class structure of the functional or local odor set (compared to Fig. 3 in the main paper).