

# Supplementary Material: Efficient hierarchical Bayesian inference for spatio-temporal regression models in neuroimaging

## Summary of the proposed algorithms and derived update rules

We proposed two algorithms in the main paper, namely full and thin Dugh, which are summarized in Algorithm 1 and Algorithm 2, respectively.

---

### Algorithm 1: Full Dugh

---

**Input:** The lead field matrix  $\mathbf{L} \in \mathbb{R}^{M \times N}$  and  $G$  trials of measurement vectors  $\{\mathbf{Y}_g\}_{g=1}^G$ , where  $\mathbf{Y}_g \in \mathbb{R}^{M \times T}$ .

**Result:** Estimates of the source and noise variances  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ , the temporal covariance  $\mathbf{B}$ , and the posterior mean  $\{\bar{\mathbf{x}}_g\}_{g=1}^G$  and covariance  $\Sigma_{\mathbf{x}}$  of the sources.

- 1 Choose a random initial value for  $\mathbf{B}$  as well as  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ , and construct  $\mathbf{H} = \text{diag}(\mathbf{h})$  and  $\mathbf{\Gamma} = \text{diag}([\gamma_1, \dots, \gamma_N]^\top)$ .
  - 2 Construct the augmented lead field  $\mathbf{\Phi} = [\mathbf{L}, \mathbf{I}_M]$ .
  - 3 Calculate the lead field  $\mathbf{D} = \mathbf{L} \otimes \mathbf{I}_T$  for vectorized sources.
  - 4 Calculate the prior spatio-temporal covariance for the sources as  $\Sigma_0 = \mathbf{\Gamma} \otimes \mathbf{B}$ .
  - 5 Calculate the spatial statistical covariance  $\Sigma_{\mathbf{y}} = \mathbf{\Phi} \mathbf{H} \mathbf{\Phi}^\top$ .
  - 6 Calculate the spatio-temporal statistical covariance  $\tilde{\Sigma}_{\mathbf{y}} = \Sigma_{\mathbf{y}} \otimes \mathbf{B}$ .
  - 7 Initialize  $k \leftarrow 1$
  - repeat**
  - 8     Calculate the posterior mean as  $\bar{\mathbf{x}}_g = \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_{\mathbf{y}}^{-1} \mathbf{y}_g$ , for  $g = 1, \dots, G$ , where  $\mathbf{y}_g = \text{vec}(\mathbf{Y}_g^\top) \in \mathbb{R}^{MT \times 1}$ .
  - 9     Calculate  $\mathbf{M}_{\text{time}}^k$  based on Eq. (6), and update  $\mathbf{B}$  based on Eq. (7) according to the Riemannian update on the manifold of P.D. matrices.
  - 10    Calculate  $\mathbf{M}_{\text{SN}}^k$  based on Eq. (9), and update  $\mathbf{H}$  based on Eq. (10).
  - 11     $k \leftarrow k + 1$
  - until** stopping condition is satisfied:  $\|\bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k\|_2^2 \leq \epsilon$  or  $k = k_{\text{max}}$ ;
  - 12 Calculate the posterior covariance as  $\Sigma_{\mathbf{x}} = \Sigma_0 - \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_{\mathbf{y}}^{-1} \mathbf{D} \Sigma_0$ .
-

---

**Algorithm 2:** Thin Dugh

---

**Input :** The lead field matrix  $\mathbf{L} \in \mathbb{R}^{M \times N}$ , and  $G$  trials of measurement vectors  $\{\mathbf{Y}_g\}_{g=1}^G$ , where  $\mathbf{Y}_g \in \mathbb{R}^{M \times T}$ .

**Result:** Estimates of the source and noise variances  $\mathbf{h} = [\gamma_1, \dots, \gamma_N, \sigma_1^2, \dots, \sigma_M^2]^\top$ , the temporal covariance  $\mathbf{B}$ , and the posterior mean  $\{\bar{\mathbf{x}}_g\}_{g=1}^G$ .

- 1 Choose a random initial value for  $\mathbf{p}$  as well as  $\mathbf{h}$ , and construct  $\mathbf{H} = \text{diag}(\mathbf{h})$  and  $\mathbf{P} = \text{diag}(\mathbf{p})$ .
  - 2 Construct  $\mathbf{B} = \mathbf{Q}\mathbf{P}\mathbf{Q}^H$ , where  $\mathbf{Q} = [\mathbf{I}_M, \mathbf{0}]\mathbf{F}_L$  with  $L = 2T + 1$  and  $\mathbf{F}_L$  as DFT.
  - 3 Construct the augmented lead field  $\Phi := [\mathbf{L}, \mathbf{I}_M]$ .
  - 4 Calculate the lead field  $\mathbf{D} = \mathbf{L} \otimes \mathbf{I}_T$  for vectorized sources.
  - 5 Calculate the statistical covariance  $\Sigma_{\mathbf{y}} = \Phi\mathbf{H}\Phi^\top$ .
  - 6 Calculate the statistical covariance  $\Sigma_{\mathbf{y}} = \Phi\mathbf{H}\Phi^\top$ .
  - 7 Calculate the spatio-temporal statistical covariance  $\tilde{\Sigma}_{\mathbf{y}} = \Sigma_{\mathbf{y}} \otimes \mathbf{B}$ .
  - 8 Initialize  $k \leftarrow 1$
  - repeat**
  - 9     Calculate the posterior mean efficiently based on Eq. (19) as  
       $\bar{\mathbf{x}}_g = \text{tr}(\mathbf{Q}\mathbf{P}(\mathbf{\Pi} \odot \mathbf{Q}^H \mathbf{Y}_g^\top \mathbf{U}_{\mathbf{x}})(\mathbf{U}_{\mathbf{x}}^\top \mathbf{L}\mathbf{\Gamma}^\top))$ , where  $\mathbf{L}\mathbf{\Gamma}\mathbf{L}^\top = \mathbf{U}_{\mathbf{x}}\mathbf{D}_{\mathbf{x}}\mathbf{U}_{\mathbf{x}}^\top$  and  $[\mathbf{\Pi}]_{l,m} = \frac{1}{\sigma_m^2 + p_l d_m}$  for  $l = 1, \dots, L$  and  $m = 1, \dots, M$ .
  - 10    Calculate  $\mathbf{M}_{\text{time}}^k$  based on Eq. (6), and update  $\mathbf{B}$  based on Eq. (16) according to Riemannian update for Toeplitz matrices using circulant embedding.
  - 11    Calculate  $\mathbf{M}_{\text{SN}}^k$  based on Eq. (9), and update  $\mathbf{H}$  based on Eq. (10).
  - 12     $k \leftarrow k + 1$
  - until** *stopping condition is satisfied:*  $\|\bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k\|_2^2 \leq \epsilon$  or  $k = k_{\text{max}}$ ;
  - 13 Calculate the posterior covariance as  $\Sigma_{\mathbf{x}} = \Sigma_0 - \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_{\mathbf{y}}^{-1} \mathbf{D} \Sigma_0$ .
- 

## A Derivation of Type-II Bayesian cost function for full-structural spatio-temporal models

In this section, we provide a detailed derivation of Type-II Bayesian learning for full-structural spatio-temporal models. To this end, we first briefly explain the *multiple measurement vector* (MMV) model and then formulate Type-II Bayesian learning with full-structural spatio-temporal covariance structure for this setting. Note that, to simplify the problem, we first present the derivations of the MMV model only for a single trial. We later extend this simplified setting to the multi-trials case.

### A.1 Multiple measurement vector (MMV) model

In M/EEG brain source imaging, a sequence of measurement vectors are often available. Thus, the following *multiple measurement vector* (MMV) model can be formulated1:

$$\mathbf{Y} = \mathbf{L}\mathbf{X} + \mathbf{E},$$

where  $\mathbf{Y} = [\mathbf{y}(1), \dots, \mathbf{y}(T)] \in \mathbb{R}^{M \times T}$  consists of  $T$  measurement vectors for a sequence of  $T$  time samples.  $\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(T)] \in \mathbb{R}^{N \times T}$  is the desired solution matrix (the amplitude of  $N$  brain sources during  $T$  time samples in our setting), and  $\mathbf{E}$  is an unknown noise matrix. A key assumption in the MMV model is that the support (i.e., the indices of the nonzero entries) of every column in  $\mathbf{X}$  is identical (referred to as the *common sparsity assumption* in the literature). The number of nonzero rows in  $\mathbf{X}$  needs to be below a threshold to ensure unique and global solution. This implies that  $\mathbf{X}$  has only a small number of non-zero rows. It has been shown that the recovery of the support can be greatly improved by increasing the number of measurements [50–52].

### A.2 Type-II Bayesian cost function for full-structural spatio-temporal models

To exploit temporal correlations between measurements, we first assume that the voxels are mutually independent. Given the column vector  $\gamma = [\gamma_1, \dots, \gamma_N]^\top$  and a Gaussian probability density for each brain source, the prior distribution with time correlation is modeled as follows:

$$p(\mathbf{X}_i | \gamma_i, \mathbf{B}) \sim \mathcal{N}(0, \gamma_i \mathbf{B}), \quad i = 1, \dots, N. \quad (20)$$

$\mathbf{X}_i$  denotes the  $i$ -th row of source matrix  $\mathbf{X}$  and models the probability distribution of the  $i$ -th brain source. Note that  $\gamma_i$  is a non-negative hyper-parameter that controls the row sparsity of  $\mathbf{X}$ ; i.e., the values of source  $\mathbf{X}_i$  become all zero if  $\gamma_i = 0$ . Finally,  $\mathbf{B}$  is a positive definite matrix that captures the time correlation structure, which is assumed to be shared across all sources. The goal is to obtain the prior distribution of sources,  $p(\mathbf{X}|\gamma, \mathbf{B})$ , by estimating the hyper-parameters,  $\{\gamma, \mathbf{B}\}$ . Next, we reformulate the joint MMV model of all sources using vectorization of matrices and Kronecker product operations:

$$\mathbf{y} = \mathbf{D}\mathbf{x} + \mathbf{e} ,$$

where  $\mathbf{y} = \text{vec}(\mathbf{Y}^\top) \in \mathbb{R}^{MT \times 1}$ ,  $\mathbf{x} = \text{vec}(\mathbf{X}^\top) \in \mathbb{R}^{NT \times 1}$ ,  $\mathbf{e} = \text{vec}(\mathbf{E}^\top) \in \mathbb{R}^{MT \times 1}$  and  $\mathbf{D} = \mathbf{L} \otimes \mathbf{I}_T$ .

The prior distribution of  $\mathbf{x}$  is given as

$$p(\mathbf{x}|\gamma, \mathbf{B}) \sim \mathcal{N}(0, \Sigma_0)$$

where  $\Sigma_0$  is defined as

$$\Sigma_0 = \begin{bmatrix} \gamma_1 \mathbf{B} & & \\ & \ddots & \\ & & \gamma_N \mathbf{B} \end{bmatrix} = \Gamma \otimes \mathbf{B} ,$$

in which  $\Gamma = \text{diag}(\gamma) = \text{diag}(\gamma_1, \dots, \gamma_N)$ .

Similarly, we may assume zero-mean Gaussian noise with covariance  $\Sigma_e = \Lambda \otimes \Upsilon$ , where  $\mathbf{e} \sim \mathcal{N}(0, \Sigma_e)$ , and  $\Lambda$  and  $\Upsilon$  denote the spatial and temporal noise covariance matrices, respectively. Here, we use the same prior for the temporal structure of noise and sources, i.e.,  $\Upsilon = \mathbf{B}$ .

The parameters of the spatio-temporal Type-II model are the unknown source, noise and temporal covariance matrices, i.e.,  $\Gamma$ ,  $\Lambda$ , and  $\mathbf{B}$ . The unknown parameters  $\Gamma$ ,  $\Lambda$ , and  $\mathbf{B}$  are optimized based on the current estimates of the source, noise and temporal covariances in an alternating iterative process. Given initial estimates of  $\Gamma$ ,  $\Lambda$ , and  $\mathbf{B}$ , the posterior distribution of the sources is a Gaussian of the form  $p(\mathbf{x}|\mathbf{y}, \Gamma, \Lambda, \mathbf{B}) \sim \mathcal{N}(\bar{\mathbf{x}}, \Sigma_x)$ , whose mean and covariance are obtained as follows:

$$\bar{\mathbf{x}} = \Sigma_0 \mathbf{D}^\top (\Lambda \otimes \mathbf{B} + \mathbf{D} \Sigma_0 \mathbf{D}^\top)^{-1} \mathbf{y} = \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_y^{-1} \mathbf{y} , \quad (21)$$

$$\Sigma_x = \Sigma_0 - \Sigma_0 \mathbf{D}^\top \tilde{\Sigma}_y^{-1} \mathbf{D} \Sigma_0 , \quad (22)$$

where  $\Sigma_y = \mathbf{L} \Gamma \mathbf{L}^\top + \Lambda$ , and where  $\tilde{\Sigma}_y = \Sigma_y \otimes \mathbf{B}$  denotes the spatio-temporal variant of statistical model covariance matrix. The estimated posterior parameters  $\bar{\mathbf{x}}$  and  $\Sigma_x$  are then in turn used to update  $\Gamma$ ,  $\Lambda$ , and  $\mathbf{B}$  as the minimizers of the negative log of the marginal likelihood  $p(\mathbf{Y}|\Gamma, \Lambda, \mathbf{B})$ , which is given by

$$\mathcal{L}_{\text{kron}}(\Gamma, \Lambda, \mathbf{B}) = \log |\tilde{\Sigma}_y| + \text{tr} \left( \mathbf{y}^\top \tilde{\Sigma}_y^{-1} \mathbf{y} \right) . \quad (23)$$

Using the same temporal covariance prior for noise and sources, i.e.,  $\Upsilon = \mathbf{B}$ , the statistical model covariance matrix,  $\tilde{\Sigma}_y$ , can be written as:

$$\begin{aligned} \tilde{\Sigma}_y &= \Lambda \otimes \Upsilon + (\mathbf{D} \Sigma_0 \mathbf{D}^\top) = \Lambda \otimes \Upsilon + \left( (\mathbf{L} \otimes \mathbf{I}^\top) (\Gamma \otimes \mathbf{B}) (\mathbf{L} \otimes \mathbf{I}^\top)^\top \right) \\ &= \Lambda \otimes \Upsilon + (\mathbf{L} \Gamma \mathbf{L}^\top \otimes \mathbf{B}) \stackrel{(\Upsilon = \mathbf{B})}{=} (\Lambda + \mathbf{L} \Gamma \mathbf{L}^\top) \otimes \mathbf{B} \\ &= \Sigma_y \otimes \mathbf{B} , . \end{aligned} \quad (24)$$

which leads to the following spatio-temporal Type-II Bayesian learning cost function:

$$\begin{aligned} \mathcal{L}_{\text{kron}}(\Gamma, \Lambda, \mathbf{B}) &= \log |\tilde{\Sigma}_y| + \text{tr} \left( \mathbf{y}^\top \tilde{\Sigma}_y^{-1} \mathbf{y} \right) \\ &= \log |\Sigma_y \otimes \mathbf{B}| + \text{tr} \left( \mathbf{y}^\top (\Sigma_y \otimes \mathbf{B})^{-1} \mathbf{y} \right) \\ &= \log (|\Sigma_y|^T |\mathbf{B}|^M) + \text{tr} \left( \mathbf{y}^\top (\Sigma_y \otimes \mathbf{B})^{-1} \mathbf{y} \right) . \end{aligned} \quad (25)$$

Here, we assume the presence of  $G$  sample blocks  $\mathbf{Y}_g \in \mathbb{R}^{M \times T}$ , for  $g = 1, \dots, G$ . These block samples can be obtained by segmenting a time series into smaller parts that are assumed to independent

and identically distributed. These blocks may represent epochs, trials or experimental tasks depending on the applications.  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  can then be reformulated as

$$\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_{\mathbf{y}}| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^{\top}) \quad (26)$$

by applying the following matrix equality to Eq. (25):

$$\text{tr}(\mathbf{A}^{\top} \mathbf{B} \mathbf{C} \mathbf{D}^{\top}) = \text{vec}(\mathbf{A})^{\top} (\mathbf{D} \otimes \mathbf{B}) \text{vec}(\mathbf{C}) .$$

## B Proof of Theorem 1

*Proof.* We start by recalling  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  in Eq. (4):

$$\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_{\mathbf{y}}| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^{\top}).$$

Let  $\mathbf{\Sigma}_{\mathbf{y}}^k$ ,  $\mathbf{\Gamma}^k$ , and  $\mathbf{\Lambda}^k$  be the values of statistical model covariance and the source and noise covariances at the  $k$ -th iteration, respectively. By ignoring terms that do not depend on  $\mathbf{B}$ ,  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  can be written as follows:

$$\begin{aligned} \mathcal{L}_{\text{kron}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) &= M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr} \left( (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^{\top} \right) \\ &= \log |\mathbf{B}| + \frac{1}{MG} \sum_{g=1}^G \text{tr} \left( (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^{\top} \right) \\ &= \log |\mathbf{B}| + \text{tr} \left( \mathbf{B}^{-1} \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^{\top} (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{Y}_g \right) \\ &= \log |\mathbf{B}| + \text{tr} (\mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k) , \end{aligned} \quad (27)$$

where  $\mathbf{M}_{\text{time}}^k := \frac{1}{MG} \sum_{g=1}^G \mathbf{Y}_g^{\top} (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{Y}_g$ .

By virtue of the concavity of the log-determinant function and its first order Taylor expansion around  $\mathbf{B}^k$ , the following inequality holds:

$$\begin{aligned} \mathcal{L}_{\text{kron}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) &= \log |\mathbf{B}| + \text{tr} (\mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k) \\ &\leq \log |\mathbf{B}^k| + \text{tr} \left( (\mathbf{B}^k)^{-1} (\mathbf{B} - \mathbf{B}^k) \right) + \text{tr} (\mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k) \\ &= \log |\mathbf{B}^k| + \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B} \right) - \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B}^k \right) + \text{tr} (\mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k) \\ &= \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B} \right) + \text{tr} (\mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k) + \text{const} \\ &= \mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) + \text{const} . \end{aligned} \quad (28)$$

Note that constant values in (28) do not depend on  $\mathbf{B}$ ; hence, they can be ignored in the optimization procedure. Hence, we have shown that minimizing Eq. (4) with respect to  $\mathbf{B}$  is equivalent to minimizing  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$ , which concludes the proof.  $\square$

## C Proof of Theorem 2

Before presenting the proof, the subsequent definitions and propositions are required:

**Definition 1** (Geodesic path). *Let  $\mathcal{M}$  be a Riemannian manifold, i.e., a differentiable manifold whose tangent space is endowed with an inner product that defines local Euclidean structures. Then, a geodesic between two points on  $\mathcal{M}$ , denoted by  $\mathbf{p}_0, \mathbf{p}_1 \in \mathcal{M}$ , is defined as the shortest connecting path between those two points along the manifold,  $\zeta_l(\mathbf{p}_0, \mathbf{p}_1) \in \mathcal{M}$  for  $l \in [0, 1]$ , where  $l = 0$  and  $l = 1$  defines the starting and end points of the path, respectively.*

In the current context,  $\zeta_l(\mathbf{p}_0, \mathbf{p}_1)$  defines a geodesic curve on the P.D. manifold joining two P.D. matrices,  $\mathbf{P}_0, \mathbf{P}_1 > 0$ . The specific pair of matrices we will deal with is  $\{\mathbf{B}^k, \mathbf{M}_{\text{time}}^k\}$ .

**Definition 2** (Geodesic on the P.D. manifold), *Geodesics on the manifold of P.D. matrices can be shown to form a cone within the embedding space. We denote this manifold by  $\mathcal{S}_{++}$ . Assume two P.D. matrices  $\mathbf{P}_0, \mathbf{P}_1 \in \mathcal{S}_{++}$ . Then, for  $l \in [0, 1]$ , the geodesic curve joining  $\mathbf{P}_0$  to  $\mathbf{P}_1$  is defined as [53, Chapter. 6]:*

$$\xi_l(\mathbf{P}_0, \mathbf{P}_1) = (\mathbf{P}_0)^{\frac{1}{2}} \left( (\mathbf{P}_0)^{-1/2} \mathbf{P}_1 (\mathbf{P}_0)^{-1/2} \right)^l (\mathbf{P}_0)^{\frac{1}{2}} \quad l \in [0, 1]. \quad (29)$$

Note that  $\mathbf{P}_0$  and  $\mathbf{P}_1$  are obtained as the starting and end points of the geodesic path by choosing  $l = 0$  and  $l = 1$ , respectively. The midpoint of the geodesic, obtained by setting  $l = \frac{1}{2}$ , is called the *geometric mean*. Note that, according to Definition 2, the following equality holds :

$$\begin{aligned} \xi_l(\mathbf{B}_0, \mathbf{B}_1)^{-1} &= \left( (\mathbf{B}_0)^{1/2} \left( (\mathbf{B}_0)^{-1/2} \mathbf{B}_1 (\mathbf{B}_0)^{-1/2} \right)^l (\mathbf{B}_0)^{1/2} \right)^{-1} \\ &= \left( (\mathbf{B}_0)^{-1/2} \left( (\mathbf{B}_0)^{1/2} (\mathbf{B}_1)^{-1} (\mathbf{B}_0)^{1/2} \right)^l (\mathbf{B}_0)^{-1/2} \right) = \xi_l(\mathbf{B}_0^{-1}, \mathbf{B}_1^{-1}). \end{aligned} \quad (30)$$

**Definition 3** (Geodesic convexity). *Let  $\mathbf{p}_0$  and  $\mathbf{p}_1$  be two arbitrary points on a subset  $\mathcal{A}$  of a Riemannian manifold  $\mathcal{M}$ . Then, a real-valued function  $f : \mathcal{A} \rightarrow \mathbb{R}$  with domain  $\mathcal{A} \subset \mathcal{M}$  is called geodesic convex (g-convex) if the following relation holds:*

$$f(\zeta_l(\mathbf{p}_0, \mathbf{p}_1)) \leq lf(\mathbf{p}_0) + (1-l)f(\mathbf{p}_1), \quad (31)$$

where  $l \in [0, 1]$  and  $\zeta(\mathbf{p}_0, \mathbf{p}_1)$  denotes the geodesic path connecting two points  $\mathbf{p}_0$  and  $\mathbf{p}_1$  as defined in Definition 1. Thus, in analogy to classical convexity, the function  $f$  is g-convex if every geodesic  $\zeta(\mathbf{p}_0, \mathbf{p}_1)$  of  $\mathcal{M}$  between  $\mathbf{p}_0, \mathbf{p}_1 \in \mathcal{A}$ , lies in the g-convex set  $\mathcal{A}$ . Note that the set  $\mathcal{A} \subset \mathcal{M}$  is called g-convex, if any geodesics joining an arbitrary pair of points lies completely in  $\mathcal{A}$ .

**Remark 2.** *Note that g-convexity is a generalization of classical (linear) convexity to non-Euclidean (non-linear) geometry and metric spaces. Therefore, it is straightforward to show that all convex functions in Euclidean geometry are also g-convex, where the geodesics between pairs of matrices are simply line segments:*

$$\zeta_l(\mathbf{p}_0, \mathbf{p}_1) = l\mathbf{p}_0 + (1-l)\mathbf{p}_1. \quad (32)$$

For the sake of brevity, we omit a detailed theoretical introduction of g-convexity, and only borrow a result from [54]. Interested readers are referred to [55, Chapter 1] for a gentle introduction to this topic, and [56, Chapter. 2]; [57–64] for more in-depth technical details. Now we are ready to state the proof, which parallels the one provided in Zadeh et al. [54, Theorem. 3].

*Proof.* We proceed in two steps. First, we consider P.D. manifolds and express (31) in terms of geodesic paths and functions that lie on this particular space. We then show that  $\mathcal{L}_{\text{conv}}^{\text{time}}(\Gamma^k, \Lambda^k, \mathbf{B})$  is strictly g-convex on this specific domain. In the second step, we then derive the update rule proposed in (7).

### C.1 Part I: G-convexity of the majorizing cost function

We consider geodesics along the P.D. manifold by setting  $\zeta_l(\mathbf{p}_0, \mathbf{p}_1)$  to  $\xi_l(\mathbf{B}_0, \mathbf{B}_1)$  as presented in Definition 2, and define  $f(\cdot)$  to be  $f(\mathbf{B}) = \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B} \right) + \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1})$ , representing the cost function  $\mathcal{L}_{\text{conv}}^{\text{time}}(\Gamma^k, \Lambda^k, \mathbf{B})$ .

We now show that  $f(\mathbf{B})$  is strictly g-convex on this specific domain. For continuous functions as considered in this paper, fulfilling (31) for  $f(\mathbf{B})$  and  $\xi_l(\mathbf{B}_0, \mathbf{B}_1)$  with  $l = 1/2$  is sufficient for strict g-convexity according to *mid-point convexity* [65]:

$$\begin{aligned} &\text{tr} \left( (\mathbf{B}^k)^{-1} \xi_{1/2}(\mathbf{B}_0, \mathbf{B}_1) \right) + \text{tr} \left( \mathbf{M}_{\text{time}}^k \xi_{1/2}(\mathbf{B}_0, \mathbf{B}_1)^{-1} \right) \\ &< \frac{1}{2} \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B}_0 \right) + \frac{1}{2} \text{tr} \left( \mathbf{M}_{\text{time}}^k \mathbf{B}_0^{-1} \right) \\ &\quad + \frac{1}{2} \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B}_1 \right) + \frac{1}{2} \text{tr} \left( \mathbf{M}_{\text{time}}^k \mathbf{B}_1^{-1} \right). \end{aligned} \quad (33)$$

Given  $(\mathbf{B}^k)^{-1} \in \mathcal{S}_{++}$ , i.e.,  $(\mathbf{B}^k)^{-1} > 0$  and the operator inequality [53, Chapter. 4]

$$\xi_{1/2}(\mathbf{B}_0, \mathbf{B}_1) \prec \frac{1}{2}\mathbf{B}_0 + \frac{1}{2}\mathbf{B}_1, \quad (34)$$

we have:

$$\text{tr} \left( (\mathbf{B}^k)^{-1} \xi_{1/2}(\mathbf{B}_0, \mathbf{B}_1) \right) < \frac{1}{2} \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B}_0 \right) + \frac{1}{2} \text{tr} \left( (\mathbf{B}^k)^{-1} \mathbf{B}_1 \right), \quad (35)$$

which is derived by multiplying both sides of Eq. (34) with  $(\mathbf{B}^k)^{-1}$  followed by taking the trace on both sides.

Similarly, we can write the operator inequality for  $\{\mathbf{B}_0^{-1}, \mathbf{B}_1^{-1}\}$  using Eq. (30) as:

$$\xi_{1/2}(\mathbf{B}_0, \mathbf{B}_1)^{-1} = \xi_{1/2}(\mathbf{B}_0^{-1}, \mathbf{B}_1^{-1}) \prec \frac{1}{2}\mathbf{B}_0^{-1} + \frac{1}{2}\mathbf{B}_1^{-1}. \quad (36)$$

Multiplying both sides of Eq. (36) by  $\mathbf{M}_{\text{time}}^k \in \mathcal{S}_{++}$  and applying the trace operator on both sides leads to:

$$\text{tr} \left( \mathbf{M}_{\text{time}}^k \xi_{1/2}(\mathbf{B}_0, \mathbf{B}_1)^{-1} \right) < \frac{1}{2} \text{tr} \left( \mathbf{M}_{\text{time}}^k \mathbf{B}_0^{-1} \right) + \frac{1}{2} \text{tr} \left( \mathbf{M}_{\text{time}}^k \mathbf{B}_1^{-1} \right). \quad (37)$$

Summing up (35) and (37) proves inequality (33) and concludes the first part of the proof.

## C.2 Part II: Derivation of the update rule in Eq. (7)

We now present the second part of the proof by deriving the update rule in Eq. (7). Since the cost function  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$  is strictly g-convex, its optimal solution in the  $k$ -th iteration is unique. More concretely, the optimum can be analytically derived by taking the derivative of Eq. (7) and setting the result to zero as follows:

$$\nabla \mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B}) = (\mathbf{B}^k)^{-1} - \mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k \mathbf{B}^{-1} = 0, \quad (38)$$

which results in

$$\mathbf{B} (\mathbf{B}^k)^{-1} \mathbf{B} = \mathbf{M}_{\text{time}}^k. \quad (39)$$

This solution is known as the *Riccati equation* and is the geometric mean between  $\mathbf{B}^k$  and  $\mathbf{M}_{\text{time}}^k$  [61, 66]:

$$\mathbf{B}^{k+1} \leftarrow (\mathbf{B}^k)^{\frac{1}{2}} \left( (\mathbf{B}^k)^{-1/2} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1/2} \right)^{\frac{1}{2}} (\mathbf{B}^k)^{\frac{1}{2}}.$$

Deriving the update rule in Eq. (7) concludes the second part of the proof of Theorem 2.  $\square$

## D Proof of Theorem 3

*Proof.* Analogous to the proof of Theorem 1 in Appendix B, we start by recalling  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  in Eq. (4):

$$\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}) = T \log |\mathbf{\Sigma}_{\mathbf{y}}| + M \log |\mathbf{B}| + \frac{1}{G} \sum_{g=1}^G \text{tr}(\mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{Y}_g \mathbf{B}^{-1} \mathbf{Y}_g^{\top}).$$

Let  $\mathbf{B}^k$  be the value of the temporal covariance matrix learned using Eq. (7) at  $k$ -th iteration. Then, by ignoring the term  $M \log |\mathbf{B}^k|$  that is only a function of  $\mathbf{B}^k$ ,  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  can be written as follows:

$$\begin{aligned} \mathcal{L}_{\text{kron}}^{\text{space}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}^k) &= T \log |\mathbf{\Sigma}_{\mathbf{y}}| + \frac{1}{G} \sum_{g=1}^G \text{tr} \left( \mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{Y}_g (\mathbf{B}^k)^{-1} \mathbf{Y}_g^{\top} \right) \\ &= \log |\mathbf{\Sigma}_{\mathbf{y}}| + \frac{1}{TG} \sum_{g=1}^G \text{tr} \left( \mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{Y}_g (\mathbf{B}^k)^{-1} \mathbf{Y}_g^{\top} \right) \\ &= \log |\mathbf{\Sigma}_{\mathbf{y}}| + \text{tr} \left( \mathbf{\Sigma}_{\mathbf{y}}^{-1} \frac{1}{TG} \sum_{g=1}^G \mathbf{Y}_g (\mathbf{B}^k)^{-1} \mathbf{Y}_g^{\top} \right) \\ &= \log |\mathbf{\Sigma}_{\mathbf{y}}| + \text{tr} \left( \mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k \right), \end{aligned} \quad (40)$$

where  $\mathbf{M}_{\text{space}}^k := \frac{1}{TG} \sum_{g=1}^G \mathbf{Y}_g (\mathbf{B}^k)^{-1} \mathbf{Y}_g^\top$ .

Similar to the argument made in Appendix B, a first order Taylor expansion of the log-determinant function around  $\Sigma_{\mathbf{y}}$  provides the following inequality:

$$\begin{aligned} \log |\Sigma_{\mathbf{y}}| &\leq \log |\Sigma_{\mathbf{y}}^k| + \text{tr} \left( (\Sigma_{\mathbf{y}}^k)^{-1} (\Sigma_{\mathbf{y}} - \Sigma_{\mathbf{y}}^k) \right) \\ &= \log |\Sigma_{\mathbf{y}}^k| + \text{tr} \left( (\Sigma_{\mathbf{y}}^k)^{-1} \Sigma_{\mathbf{y}} \right) - \text{tr} \left( (\Sigma_{\mathbf{y}}^k)^{-1} \Sigma_{\mathbf{y}}^k \right) \\ &= \text{tr}(\Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}) + \text{const} , \end{aligned} \quad (41)$$

where the last step is derived using the augmented source and noise covariances,  $\mathbf{H} := [\Gamma, \mathbf{0}; \mathbf{0}, \Lambda]$ ,  $\Phi := [\mathbf{L}, \mathbf{I}]$  and  $\Sigma_{\mathbf{y}} = \Phi \mathbf{H} \Phi^\top$ .

By inserting Eq. (40) into Eq. (41), the first term of Eq. (8),  $\text{tr}(\Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H})$ , can be directly inferred:

$$\begin{aligned} \mathcal{L}_{\text{kron}}^{\text{space}}(\Gamma, \Lambda, \mathbf{B}^k) &= \mathcal{L}_{\text{kron}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k) = \log |\Sigma_{\mathbf{y}}| + \text{tr}(\Sigma_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k) \\ &\leq \text{tr}(\Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}) + \text{tr}(\Sigma_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k) + \text{const} , \end{aligned} \quad (42)$$

We further show how the second term in Eq. (8) can be derived. To this end, we construct an upper bound on  $\text{tr}(\Sigma_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k)$  using an inequality derived from the Schur complement of  $\Sigma_{\mathbf{y}}$ . Before presenting this inequality, the subsequent definition of the Schur complement of matrix  $\Sigma_{\mathbf{y}}$  is required:

**Definition 4.** For a positive semidefinite (PSD) matrix  $\Sigma_{\mathbf{y}}$ , and a partitioning

$$\mathbf{X} = \begin{bmatrix} \mathbf{D} & \mathbf{G} \\ \mathbf{G}^\top & \mathbf{B} \end{bmatrix} , \quad (43)$$

its Schur complement is defined as

$$\mathbf{S} := \mathbf{D} - \mathbf{G} \Sigma_{\mathbf{y}}^{-1} \mathbf{G}^\top \quad (44)$$

$$(45)$$

The Schur complement condition states that the matrix  $\mathbf{X}$  is PSD,  $\mathbf{X} \geq \mathbf{0}$ , if and only if the Schur complement of  $\Sigma_{\mathbf{y}}$  is PSD,  $\mathbf{S} \geq \mathbf{0}$ .

Now we are ready to construct an upper bound on  $\text{tr}(\Sigma_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k)$ . To this end, we show that  $\text{tr}(\Sigma_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k)$  can be majorized as follows:

$$\text{tr}(\Sigma_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k) \leq \text{tr}(\mathbf{H}^k \Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} \mathbf{M}_{\text{space}}^k (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-1}) . \quad (46)$$

By defining  $\mathbf{V}$  as:

$$\mathbf{V} = \begin{bmatrix} (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-\frac{1}{2}} \\ \Phi \mathbf{H}^{\frac{1}{2}} \end{bmatrix} , \quad (47)$$

the PSD property of  $\mathbf{S}$  can be inferred as:

$$\mathbf{S} = \begin{bmatrix} (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-1} \mathbf{H}^k \Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} & \mathbf{I} \\ \mathbf{I} & \Phi \mathbf{H} \Phi^\top \end{bmatrix} = \mathbf{V} \mathbf{V}^\top \geq \mathbf{0} . \quad (48)$$

By employing the definition of the Schur complement with  $\mathbf{D} = (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-1} \mathbf{H}^k \Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1}$ ,  $\mathbf{G} = \mathbf{I}$  and  $\Sigma_{\mathbf{y}} = \Phi \mathbf{H} \Phi^\top$ , we have:

$$(\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-1} \mathbf{H}^k \Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} \geq (\Phi \mathbf{H} \Phi^\top)^{-1} . \quad (49)$$

The inequality in Eq. (46) can be directly inferred by multiplying  $\mathbf{M}_{\text{space}}^k$  to both sides of Eq. (49), applying trace operator, and rearranging the arguments in the trace operator:

$$\begin{aligned} \text{tr}(\mathbf{M}_{\text{space}}^k \Sigma_{\mathbf{y}}^{-1}) &\leq \text{tr}(\mathbf{M}_{\text{space}}^k (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-1} \mathbf{H}^k \Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1}) \\ &= \text{tr}(\mathbf{H}^k \Phi^\top (\Sigma_{\mathbf{y}}^k)^{-1} \mathbf{M}_{\text{space}}^k (\Sigma_{\mathbf{y}}^k)^{-1} \Phi \mathbf{H}^k \mathbf{H}^{-1}) \\ &= \text{tr}(\mathbf{M}_{\text{SN}}^k \mathbf{H}^{-1}) . \end{aligned} \quad (50)$$

By inserting Eq. (50) into Eq. (42), we have

$$\begin{aligned}\mathcal{L}_{\text{kron}}^{\text{space}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B}^k) &= \mathcal{L}_{\text{kron}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k) \leq \text{tr}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \mathbf{H}) + \text{tr}(\mathbf{\Sigma}_{\mathbf{y}}^{-1} \mathbf{M}_{\text{space}}^k) + \text{const} \\ &\leq \text{tr}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \mathbf{H}) + \text{tr}(\mathbf{M}_{\text{SN}}^k \mathbf{H}^{-1}) + \text{const} \\ &= \mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k) + \text{const} .\end{aligned}\quad (51)$$

Note that constant values in (51) do not depend on  $\mathbf{H}$ ; hence, they can be ignored in the optimization procedure. We have shown that minimizing Eq. (4) with respect to  $\mathbf{H}$  is equivalent to minimizing  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$ , which concludes the proof.  $\square$

## E Proof of Theorem 4

*Proof.* We proceed in two steps. First, we show that  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$  is convex in  $\mathbf{h}$ . Then, we derive the update rule proposed in Eq. (10).

### E.1 Part I: Convexity of the majorizing cost function

We start the proof by constraining  $\mathbf{H}$  to the set of diagonal matrices with non-negative entries  $\mathcal{S}$ , i.e.,  $\mathcal{S} = \{\mathbf{H} \mid \mathbf{H} = \text{diag}(\mathbf{h}) = \text{diag}([h_1, \dots, h_{N+M}]^\top), h_n \geq 0, \text{ for } i = 1, \dots, N + M\}$ . We continue by reformulating the constrained optimization with respect to the source covariance matrix,

$$\mathbf{H}^{k+1} = \arg \min_{\mathbf{H} \in \mathcal{S}, \mathbf{B} = \mathbf{B}^k} \text{tr}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \mathbf{H}) + \text{tr}(\mathbf{M}_{\text{SN}}^k \mathbf{H}^{-1}), \quad (52)$$

as follows:

$$\mathbf{h}^{k+1} = \arg \min_{\mathbf{h} \geq 0, \mathbf{B} = \mathbf{B}^k} \underbrace{\text{diag}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi}) \mathbf{h} + \text{diag}(\mathbf{M}_{\text{SN}}^k) \mathbf{h}^{-1}}_{\mathcal{L}_{\text{diag}}^{\text{space}}(\mathbf{h}|\mathbf{h}^k)}, \quad (53)$$

where  $\mathbf{h}^{-1} = [h_1^{-1}, \dots, h_{N+M}^{-1}]^\top$  is defined as the element-wise inversion of  $\mathbf{h}$ . Let  $\mathbf{V}^k := \mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi}$ . Then, we rewrite  $\mathcal{L}_{\text{diag}}^{\text{space}}(\mathbf{h}|\mathbf{h}^k)$  as

$$\mathcal{L}_{\text{diag}}^{\text{space}}(\mathbf{h}|\mathbf{h}^k) = \text{diag}(\mathbf{V}^k) \mathbf{h} + \text{diag}(\mathbf{M}_{\text{SN}}^k) \mathbf{h}^{-1}. \quad (54)$$

The convexity of  $\mathcal{L}_{\text{diag}}^{\text{space}}(\mathbf{h}|\mathbf{h}^k)$  can be directly inferred from the convexity of  $\text{diag}[\mathbf{V}^k] \mathbf{h}$  and  $\text{diag}[\mathbf{M}_{\text{SN}}^k] \mathbf{h}^{-1}$  with respect to  $\mathbf{h}$  [67, Chapter. 3].

### E.2 Part II: Derivation of the update rule in Eq. (10)

We now present the second part of the proof by deriving the update rule in Eq. (10). Since the cost function  $\mathcal{L}_{\text{diag}}^{\text{space}}(\mathbf{h}|\mathbf{h}^k)$  is convex, its optimal solution in the  $k$ -th iteration is unique. Therefore, the optimization with respect to heteroscedastic source and noise variances is carried out by taking the derivative of (53) with respect to  $h_i$ , for  $n = 1, \dots, M + N$ , and setting it to zero:

$$\begin{aligned}\frac{\partial}{\partial h_i} \left( \left[ \mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \right] h_i + \left[ \mathbf{M}_{\text{SN}}^k \right] h_i^{-1} \right) \\ = \left[ \mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \right]_{i,i} - \frac{1}{(h_i)^2} \left[ \mathbf{M}_{\text{SN}}^k \right]_{i,i} \\ = 0 \quad \text{for } i = 1, \dots, N + M ,\end{aligned}$$

where  $\mathbf{\Phi}_i$  denotes the  $n$ -th column of the augmented lead field matrix. This yields the following update rule:

$$\mathbf{H}^{k+1} = \text{diag}(\mathbf{h}^{k+1}), \text{ where, } h_i^{k+1} \leftarrow \sqrt{\frac{\left[ \mathbf{M}_{\text{SN}}^k \right]_{i,i}}{\left[ \mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \right]_{i,i}}} = \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{\eta}_n^k(t))^2}{\mathbf{\Phi}_n^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi}_i}} \quad (55)$$

for  $i = 1, \dots, N + M$ .

The updates rule in Eq. (10) can be directly inferred by defining  $\mathbf{g} := \text{diag}(\mathbf{M}_{\text{SN}}^k)$  and  $\mathbf{z} := \text{diag}(\mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi})$ , which leads to:  $g_i^k = \left[ \mathbf{M}_{\text{SN}}^k \right]_{i,i}$  and  $z_i^k = \left[ \mathbf{\Phi}^\top (\mathbf{\Sigma}_{\mathbf{y}}^k)^{-1} \mathbf{\Phi} \right]_{i,i}$ . This concludes the proof.  $\square$



## F Champagne with heteroscedastic noise learning

Interestingly, identical update rules as those proposed in Champagne [33] and heteroscedastic noise learning [68] can be derived for source and noise variances, respectively, by selecting the corresponding indices of matrix  $\mathbf{H}$  associated to noise and source covariances.

### F.1 Update rule for source variances

Given  $[\Phi]_{1:M,1:N} = \mathbf{L}$ ,  $[\mathbf{H}]_{1:N,1:N} = \mathbf{\Gamma}$ , and  $[\bar{\eta}(t)]_{1:N} = \bar{\mathbf{x}}(t)$ , the update rule for  $\mathbf{\Gamma}^{k+1} = \text{diag}(\gamma^{k+1})$  is derived by replacing  $\mathbf{H}$ ,  $\Phi$  and  $\bar{\eta}_n^k(t)$  in Eq. (55) with  $\mathbf{\Gamma}$ ,  $\mathbf{L}$  and  $\bar{\mathbf{x}}_n^k(t)$ , respectively, and defining the counterpart of  $\mathbf{M}_{\text{SN}}^k$  for sources accordingly as  $\mathbf{M}_{\text{S}}^k := \omega_{\text{S}}^k \mathbf{M}_{\text{space}}^k (\omega_{\text{S}}^k)^\top$ , where  $\omega_{\text{S}}^k := \mathbf{\Gamma}^k \mathbf{L}^\top (\Sigma_{\mathbf{y}}^k)^{-1}$ . The update rule for the source variances is then obtained as follows:

$$\gamma_n^{k+1} \leftarrow \sqrt{\frac{[\mathbf{M}_{\text{S}}^k]_{n,n}}{[\mathbf{L}^\top (\Sigma_{\mathbf{y}}^k)^{-1} \mathbf{L}]_{n,n}}} = \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}_n^k(t))^2}{\mathbf{L}_n^\top (\Sigma_{\mathbf{y}}^k)^{-1} \mathbf{L}_n}} \quad \text{for } n = 1, \dots, N, \quad (56)$$

where  $\mathbf{L}_n$  denotes the  $n$ -th column of the lead field matrix.

### F.2 Update rule for noise variances

Similarly, given  $[\Phi]_{1:M,N+1:N+M} = \mathbf{I}$ ,  $[\mathbf{H}]_{N+1:N+M,N+1:N+M} = \mathbf{\Lambda}$ , and  $[\bar{\eta}(t)]_{N+1:N+M} = \bar{\mathbf{e}}(t) := \mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}(t)$ , the update rule for  $\mathbf{\Lambda}^{k+1} = \text{diag}(\lambda^{k+1})$  is derived by replacing  $\mathbf{H}$ ,  $\Phi$  and  $\bar{\eta}_n^k(t)$  in Eq. (55) with  $\mathbf{\Lambda}$ ,  $\mathbf{I}$  and  $\bar{\mathbf{e}}_n^k(t)$ , respectively, and defining the counterpart of  $\mathbf{M}_{\text{SN}}^k$  for the noise accordingly as  $\mathbf{M}_{\text{N}}^k := \omega_{\text{N}}^k \mathbf{M}_{\text{space}}^k (\omega_{\text{N}}^k)^\top$  with  $\omega_{\text{N}}^k = \mathbf{\Lambda}^k (\Sigma_{\mathbf{y}}^k)^{-1}$ . The update rule for the noise variances is then derived as follows:

$$\lambda_m^{k+1} \leftarrow \sqrt{\frac{[\mathbf{M}_{\text{N}}^k]_{m,m}}{[(\Sigma_{\mathbf{y}}^k)^{-1}]_{m,m}}} = \sqrt{\frac{\sum_{t=1}^T (\bar{\mathbf{e}}_m^k(t))^2}{[(\Sigma_{\mathbf{y}}^k)^{-1}]_{m,m}}} \quad \text{for } m = 1, \dots, M, \quad (57)$$

which is identical to the update rule of the Champagne with heteroscedastic noise learning as presented in Cai et al. [68].

## G Proof of Theorem 5

We prove Theorem 5 by showing that the alternating update rules for  $\mathbf{B}$  and  $\mathbf{H}$ , Eqs. (7) and (10), are guaranteed to converge to a stationary point of the Bayesian Type-II likelihood  $\mathcal{L}_{\text{kron}}(\mathbf{\Gamma}, \mathbf{\Lambda}, \mathbf{B})$  Eq. (4). More generally, we prove that full Dugh is an instance of the general class of majorization-minimization (MM) algorithms, for which this property follows by construction. To this end, we first briefly review theoretical concepts behind the majorization-minimization (MM) algorithmic framework [69–72].

### G.1 Required conditions for majorization-minimization algorithms

MM is a versatile framework for optimizing general non-linear optimization programs. The main idea behind MM is to replace the original cost function in each iteration by an upper bound, also known as majorizing function, whose minimum is easy to find. Compared to other popular optimization paradigms such as (quasi)-Newton methods, MM algorithms enjoy guaranteed convergence to a stationary point [27]. The MM class covers a broad range of common optimization algorithms such as *convex-concave procedures (CCCP)* and *proximal methods* [27, Section IV], [73–75]. Such algorithms have been applied in various domains such as non-negative matrix factorization [76], graph learning [77], robust portfolio optimization in finance [78], direction of arrival (DoA) and channel estimation in wireless communications [79–82], internet of things (IoT) [83, 84], and brain

source imaging [26, 68, 85–87]. Interested readers are referred to Sun et al. [27] for an extensive list of applications on MM.

We define an original optimization problem with the objective of minimizing a continuous function  $f(\mathbf{u})$  within a closed convex set  $\mathcal{U} \subset \mathbb{R}^n$ :

$$\min_{\mathbf{u}} f(\mathbf{u}) \quad \text{subject to } \mathbf{u} \in \mathcal{U} . \quad (58)$$

Then, the idea of MM can be summarized as follows. First, construct a continuous *surrogate function*  $g(\mathbf{u}|\mathbf{u}^k)$  that *majorizes*, or upper-bounds, the original function  $f(\mathbf{u})$  and coincides with  $f(\mathbf{u})$  at a given point  $\mathbf{u}^k$ :

$$\begin{aligned} \text{[A1]} \quad & g(\mathbf{u}^k|\mathbf{u}^k) = f(\mathbf{u}^k) && \forall \mathbf{u}^k \in \mathcal{U} \\ \text{[A2]} \quad & g(\mathbf{u}|\mathbf{u}^k) \geq f(\mathbf{u}) && \forall \mathbf{u}, \mathbf{u}^k \in \mathcal{U} . \end{aligned}$$

Second, starting from an initial value  $\mathbf{u}^0$ , generate a sequence of feasible points  $\mathbf{u}^1, \mathbf{u}^2, \dots, \mathbf{u}^k, \mathbf{u}^{k+1}$  as solutions of a series of successive simple optimization problems, where

$$\text{[A3]} \quad \mathbf{u}^{k+1} := \arg \min_{\mathbf{u} \in \mathcal{U}} g(\mathbf{u}|\mathbf{u}^k) .$$

**Definition 5.** Any algorithm fulfilling conditions [A1]–[A3] is called a *Majorization Minimization (MM) algorithm*.

If a surrogate function fulfills conditions [A1]–[A3], then the value of the cost function  $f$  decreases in each iteration:

**Corollary 1.** An MM algorithm has a *descending trend property*, whereby the value of the cost function  $f$  decreases in each iteration:  $f(\mathbf{u}^{k+1}) \leq f(\mathbf{u}^k)$ .

*Proof.* To verify the descending trend in the MM framework, it is sufficient to show that  $f(\mathbf{u}^{k+1}) \leq f(\mathbf{u}^k)$ . To this end, we have  $f(\mathbf{u}^{k+1}) \leq g(\mathbf{u}^{k+1}|\mathbf{u}^k)$  from condition [A2]. Condition [A3] further states that  $g(\mathbf{u}^{k+1}|\mathbf{u}^k) \leq g(\mathbf{u}^k|\mathbf{u}^k)$ , while  $g(\mathbf{u}^k|\mathbf{u}^k) = f(\mathbf{u}^k)$  holds according to [A1]. Putting everything together, we have:

$$f(\mathbf{u}^{k+1}) \stackrel{\text{[A2]}}{\leq} g(\mathbf{u}^{k+1}|\mathbf{u}^k) \stackrel{\text{[A3]}}{\leq} g(\mathbf{u}^k|\mathbf{u}^k) \stackrel{\text{[A1]}}{=} f(\mathbf{u}^k) ,$$

which concludes the proof.  $\square$

While Corollary 1 guarantees a descending trend, convergence requires additional assumptions on particular properties of  $f$  and  $g$  [70, 71]. For the smooth functions considered in this paper, we require that the derivatives of the original and surrogate functions coincide at  $\mathbf{u}^k$ :

$$\text{[A4]} \quad \nabla g(\mathbf{u}^k|\mathbf{u}^k) = \nabla f(\mathbf{u}^k) \quad \forall \mathbf{u}^k \in \mathcal{U} .$$

We can then formulate the following, stronger, theorem:

**Theorem 8.** For an MM algorithm that additionally satisfies [A4], every limit point of the sequence of minimizers generated through [A3] is a stationary point of the original optimization problem Eq. (58).

*Proof.* A detailed proof is provided in Razaviyayn et al. [70, Theorem 1].  $\square$

Note that since we are working with smooth functions, conditions [A1]–[A4] are sufficient to prove convergence to a stationary point according to Theorem 8.

## G.2 Detailed derivation of the proof of Theorem 5

We now show that full Dugh is an instance of majorization-minimization as defined above, which fulfills Theorem 8.

*Proof.* We need to prove that conditions [A1]–[A4] are fulfilled for full Dugh. To this end, we first prove conditions [A1]–[A4] for the optimization with respect to  $\mathbf{B}$  based on the convex surrogate function in Eq. (5),  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$ . For this purpose, we recall the upper bound on  $\log |\mathbf{B}|$  in Eq. (28), which fulfills condition [A2] since it majorizes  $\log |\mathbf{B}|$  as a result of the concavity of the log-determinant function and its first-order Taylor expansion around  $\mathbf{B}^k$ . Besides, it automatically satisfies conditions [A1] and [A4] by construction, because the majorizing function in Eq. (28) is obtained through a Taylor expansion around  $\mathbf{B}^k$ . Concretely, [A1] is satisfied because the equality in Eq. (28) holds for  $\mathbf{B} = \mathbf{B}^k$ . Similarly, [A4] is satisfied because the gradient of  $\log |\mathbf{B}|$  at point  $\mathbf{B}^k$ ,  $(\mathbf{B}^k)^{-1}$  defines the linear Taylor approximation  $\log |\mathbf{B}^k| + \text{tr} \left( (\mathbf{B}^k)^{-1} (\mathbf{B} - \mathbf{B}^k) \right)$ . Thus, both gradients coincide in  $\mathbf{B}^k$  by construction. We can further prove that [A3] can be satisfied by showing that  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$  reaches its global minimum in each MM iteration. This is guaranteed if  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{\Gamma}^k, \mathbf{\Lambda}^k, \mathbf{B})$  can be shown to be convex or g-convex with respect to  $\mathbf{B}$ . To this end, we first require the subsequent proposition:

**Proposition 2.** *Any local minimum of a g-convex function over a g-convex set is a global minimum.*

*Proof.* A detailed proof is presented in Rapcsak [57, Theorem 2.1].  $\square$

Given the proof presented in Appendix C.1, we can conclude that  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{H}^k, \mathbf{B})$  is g-convex; hence, any local minimum of  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{H}^k, \mathbf{B})$  is a global minimum according to Proposition 2. This proves that condition [A3] is fulfilled and completes the proof that the optimization of Eq. (4) with respect to  $\mathbf{B}$  using the convex surrogate cost function Eq. (5) leads to an MM algorithm.

The proof of conditions [A1], [A2] and [A4] for the optimization with respect to  $\mathbf{H}$  based on the convex surrogate function in Eq. (8),  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$ , can be presented analogously. To this end, we recall the upper bound on  $\log |\Sigma_{\mathbf{y}}|$  in Eq. (41), which fulfills condition [A2] since it majorizes  $\log |\Sigma_{\mathbf{y}}|$  as a result of the concavity of the log-determinant function and its first-order Taylor expansion around  $\Sigma_{\mathbf{y}}^k$ . Besides, it automatically satisfies conditions [A1] and [A4] by construction, because the majorizing function in Eq. (41) is obtained through a Taylor expansion around  $\Sigma_{\mathbf{y}}^k$ . Concretely, [A1] is satisfied because the equality in Eq. (41) holds for  $\Sigma_{\mathbf{y}} = \Sigma_{\mathbf{y}}^k$ . Similarly, [A4] is satisfied because the gradient of  $\log |\Sigma_{\mathbf{y}}|$  at point  $\Sigma_{\mathbf{y}}^k$ ,  $(\Sigma_{\mathbf{y}}^k)^{-1}$  defines the linear Taylor approximation  $\log |\Sigma_{\mathbf{y}}^k| + \text{tr} \left[ (\Sigma_{\mathbf{y}}^k)^{-1} (\Sigma_{\mathbf{y}} - \Sigma_{\mathbf{y}}^k) \right]$ . Thus, both gradients coincide in  $\Sigma_{\mathbf{y}}^k$  by construction. We can further prove that [A3] can be satisfied by showing that  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{H}^k, \mathbf{B})$  reaches its global minimum in each MM iteration. This is guaranteed if  $\mathcal{L}_{\text{conv}}^{\text{time}}(\mathbf{H}^k, \mathbf{B})$  can be shown to be convex with respect to  $\mathbf{H} = \text{diag}(\mathbf{h})$ . Given the proof presented in Appendix E.2, we can show that [A3] is also satisfied since  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$  in Eq. (8) is a convex function with respect to  $\mathbf{h}$ . The convexity of  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$ , which ensures that condition [A3] can be satisfied using standard optimization, along with the fulfillment of conditions [A1], [A2] and [A4], ensure that Theorem 8 holds for  $\mathcal{L}_{\text{conv}}^{\text{space}}(\mathbf{H}, \mathbf{B}^k)$ . This completes the proof that the optimization of Eq. (4) with respect to  $\mathbf{H}$  using the convex surrogate cost function Eq. (8) leads to an MM algorithm that is guaranteed to converge.  $\square$

## H Proof of Proposition 1

This is a well-established classical results of the signal processing literature. Therefore, we only provide two remarks highlighting important connections between Proposition 1 and our proposed method, and refer the interested reader to Grenander and Szegö [32, Chapter 5] for a detailed proof of the theorem.

**Remark 3.** *As indicated in Babu [88], the set of embedded circulant matrices with size  $L \times L$ ,  $\mathcal{B}^L$ , does not consist exclusively of positive definite Toeplitz matrices. Therefore, we restrict ourselves to embedded circulant matrices  $\mathbf{P}$  that are also positive definite. This restriction indeed makes the diagonalization inaccurate, but we can improve the accuracy by choosing a large  $L$ . Practical evaluations have shown that choosing  $L \geq 2T - 1$  provides sufficient approximation result.*

**Remark 4.** *The Carathéodory parametrization of a Toeplitz matrix [89, Section 4.9.2] states that any PD matrix can be represented as  $\mathbf{B} = \mathbf{A}\mathbf{P}'\mathbf{A}^H$  with  $[\mathbf{A}]_{m,l} = e^{i(w_l)(m-1)}$  and  $\mathbf{P}' = \text{diag}(p'_0, p'_1, \dots, p'_{L-1})$  where  $w_l$  and  $p_l$  are specific frequencies and their corresponding*

amplitudes. By comparing the Carathéodory parametrization of  $\mathbf{B}$  with its Fourier diagonalization (Eq. (15)), it can be seen that Fourier diagonalization force the frequencies to lie on the Fourier grid, i.e.  $w_l = \frac{2\pi(l-1)}{L}$ , which indeed makes the diagonalization slightly inaccurate. The approximation accuracy can, however, be improved by increasing  $L$ . The Szegő theorem [32, 90] states that a Toeplitz matrix is asymptotically ( $L \rightarrow \infty$ ) diagonalized by the DFT matrix.

## I Proof of Theorem 6

*Proof.* We proceed in two steps. First, we show that the cost function in Eq. (13) is convex with respect to  $\mathbf{p}$ . In the second step, we then derive the update rule proposed in (16).

### I.1 Part I: Convexity of the majorizing cost function

The proof of this section parallels the one provided in [91, Proposition 4]. We start by recalling Eq. (13):

$$\mathbf{B}^* = \arg \min_{\mathbf{B} \in \mathcal{B}, \mathbf{H} = \mathbf{H}^k} \text{tr}((\mathbf{B}^k)^{-1} \mathbf{B}) + \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}). \quad (59)$$

We then show that the second term in Eq. (59) can be upper-bounded as follows:

$$\text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}) \leq \text{tr}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-1}). \quad (60)$$

By defining  $\mathbf{V}$  as

$$\mathbf{V} = \begin{bmatrix} (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-\frac{1}{2}} \\ \mathbf{Q} \mathbf{P}^{\frac{1}{2}} \end{bmatrix}, \quad (61)$$

the PSD property of  $\mathbf{S}$  can be inferred as

$$\mathbf{S} = \begin{bmatrix} (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-1} \mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} & \mathbf{I} \\ \mathbf{I} & \mathbf{Q} \mathbf{P} \mathbf{Q}^H \end{bmatrix} = \mathbf{V} \mathbf{V}^H \geq 0. \quad (62)$$

Therefore by virtue of the Schur complement with  $\mathbf{D} = (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-1} \mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1}$ ,  $\mathbf{G} = \mathbf{I}$  and  $\mathbf{B} = \mathbf{Q} \mathbf{P} \mathbf{Q}^H$ , we have:

$$(\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-1} \mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \geq (\mathbf{Q} \mathbf{P} \mathbf{Q}^H)^{-1}. \quad (63)$$

The inequality (Eq. (60)) can be directly obtained by multiplying  $\mathbf{M}_{\text{time}}^k$  to both sides of Eq. (63), applying the trace operator, using Eq. (14) and finally rearranging the terms within the trace operator:

$$\text{tr}(\mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-1} \mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1}) \geq \text{tr}(\mathbf{M}_{\text{time}}^k \mathbf{B}^{-1}). \quad (64)$$

Let  $\mathbf{B}_k = \mathbf{Q} \mathbf{P}^k \mathbf{Q}^H$  be the Fourier diagonalization of a fixed matrix  $\mathbf{B}_k$  in the  $k$ -th iteration, one can derive an efficient update rule for the temporal covariance by rewriting Eq. (13) and exploiting Propositions 1 and Eq. (60):

$$\begin{aligned} & \text{tr}((\mathbf{B}^k)^{-1} \mathbf{B}) + \text{tr}(\mathbf{B}^{-1} \mathbf{M}_{\text{time}}^k) \\ & \leq \text{tr}((\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P} \mathbf{Q}^H) + \text{tr}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k \mathbf{P}^{-1}) \\ & = \text{diag}(\mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q}) \mathbf{p} + \text{diag}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k) \mathbf{p}^{-1}, \end{aligned} \quad (65)$$

where  $\mathbf{p} = \text{vec}(\mathbf{P})$ , and  $\mathbf{p}^{-1}$  is defined as the element-wise inversion of  $\mathbf{p}$ .

We formulate the optimization problem as follows:

$$\begin{aligned} \mathcal{L}_{\text{toeplitz}}^{\text{time}}(\mathbf{p}) &= \text{diag}(\mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q}) \mathbf{p} \\ &+ \text{diag}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k) \mathbf{p}^{-1}. \end{aligned} \quad (66)$$

Let  $\mathbf{W}^k := \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q}$  and  $\mathbf{O}^k := \mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k$ . Then, we rewrite  $\mathcal{L}_{\text{toeplitz}}^{\text{time}}(\mathbf{p})$  as

$$\mathcal{L}_{\text{toeplitz}}^{\text{time}}(\mathbf{p}) = \text{diag}(\mathbf{W}^k) \mathbf{p} + \text{diag}(\mathbf{O}^k) \mathbf{p}^{-1}. \quad (67)$$

The convexity of  $\mathcal{L}_{\text{toeplitz}}^{\text{time}}(\mathbf{p})$  can be directly inferred from the convexity of  $\text{diag}[\mathbf{W}^k] \mathbf{p}$  and  $\text{diag}[\mathbf{O}^k] \mathbf{p}^{-1}$  with respect to  $\mathbf{p}$  [67, Chapter. 3].

## I.2 Part II: Derivation of the update rule in Eq. (16)

We now present the second part of the proof by deriving the update rule in Eq. (16). Since the cost function  $\mathcal{L}_{\text{toeplitz}}^{\text{time}}(\mathbf{p})$  is convex, its optimal solution in the  $k$ -th iteration is unique. More concretely, a closed-form solution of the final update rule can be obtained by taking the derivative of Eq. (67) with respect to  $\mathbf{p}$  and setting it to zero:

$$p_l^{k+1} \leftarrow \sqrt{\frac{\hat{g}_l^k}{\hat{z}_l^k}} \text{ for } l = 0, \dots, L-1, \text{ where} \quad (68)$$

$$\hat{\mathbf{g}} = \text{diag}(\mathbf{P}^k \mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{M}_{\text{time}}^k (\mathbf{B}^k)^{-1} \mathbf{Q} \mathbf{P}^k) \quad (69)$$

$$\hat{\mathbf{z}} = \text{diag}(\mathbf{Q}^H (\mathbf{B}^k)^{-1} \mathbf{Q}), \quad (70)$$

which concludes the proof.  $\square$

## J Proof of Theorem 7

*Proof.* The proof is inspired by ideas presented in Rakitsch et al. [2], Wu et al. [49], Saatçi [92] for spatio-temporal Gaussian process inference, and parallels the one proposed in Solin et al. [93]. Rakitsch et al. [2], Wu et al. [49] provide an efficient method for computing the non-convex spatio-temporal ML cost function by exploiting the compatibility between diagonalization and the Kronecker product. Here we use similar ideas to obtain the posterior mean in an efficient way.

Recalling the diagonalization of the temporal correlation matrix as  $\mathbf{B} = \mathbf{Q} \mathbf{P} \mathbf{Q}^H$  and considering the eigenvalue decomposition of  $\mathbf{L} \mathbf{\Gamma} \mathbf{L}^T$  as  $\mathbf{L} \mathbf{\Gamma} \mathbf{L}^T = \mathbf{U}_x \mathbf{D}_x \mathbf{U}_x^T$  with  $\mathbf{D}_x = \text{diag}(d_1, \dots, d_M)$ , we have:

$$\begin{aligned} \bar{\mathbf{x}}_g &= (\mathbf{\Gamma} \otimes \mathbf{B}) \mathbf{D}^T \tilde{\Sigma}_y^{-1} \mathbf{y}_g \\ &= (\mathbf{\Gamma} \otimes \mathbf{B}) (\mathbf{L} \otimes \mathbf{I})^T (\mathbf{\Lambda} \otimes \mathbf{B} + \mathbf{D} \Sigma_0 \mathbf{D}^T)^{-1} \text{vec}(\mathbf{Y}_g^T) \\ &= (\mathbf{\Gamma} \otimes \mathbf{B}) (\mathbf{L}^T \otimes \mathbf{I}) ((\mathbf{\Lambda} + \mathbf{L} \mathbf{\Gamma} \mathbf{L}^T) \otimes \mathbf{B})^{-1} \text{vec}(\mathbf{Y}_g^T) \\ &= (\mathbf{\Gamma} \mathbf{L}^T \otimes \mathbf{B}) ((\mathbf{\Lambda} + \mathbf{L} \mathbf{\Gamma} \mathbf{L}^T) \otimes \mathbf{B})^{-1} \text{vec}(\mathbf{Y}_g^T) \\ &= (\mathbf{\Gamma} \mathbf{L}^T \mathbf{U}_x \otimes \mathbf{Q} \mathbf{P}) (\mathbf{\Omega})^{-1} (\mathbf{U}_x^T \otimes \mathbf{Q}^H) \text{vec}(\mathbf{Y}_g^T) \\ &= (\mathbf{\Gamma} \mathbf{L}^T \mathbf{U}_x \otimes \mathbf{Q} \mathbf{P}) (\mathbf{\Omega})^{-1} \text{tr}(\mathbf{Q}^H \mathbf{Y}_g^T \mathbf{U}_x) \\ &= (\mathbf{\Gamma} \mathbf{L}^T \mathbf{U}_x \otimes \mathbf{Q} \mathbf{P}) \text{tr}(\mathbf{\Pi} \odot \mathbf{Q}^H \mathbf{Y}_g^T \mathbf{U}_x) \\ &= \text{tr}(\mathbf{Q} \mathbf{P} (\mathbf{\Pi} \odot \mathbf{Q}^H \mathbf{Y}_g^T \mathbf{U}_x) (\mathbf{U}_x^T \mathbf{L} \mathbf{\Gamma}^T)), \end{aligned} \quad (71)$$

where  $\odot$  denotes the Hadamard product between corresponding elements of two matrices.  $\mathbf{\Omega}$  and  $\mathbf{\Pi}$  are defined as follows:  $\mathbf{\Omega} = \mathbf{\Lambda} + \mathbf{D}_x \otimes \mathbf{P}$  and  $[\mathbf{\Pi}]_{l,m} = \frac{1}{\sigma_m^2 + p_l d_m}$  for  $l = 1, \dots, L$ ;  $m = 1, \dots, M$ . Note that the last four lines are derived based on the following matrix equality:

$$\text{tr}(\mathbf{A}^T \mathbf{B} \mathbf{C} \mathbf{D}^T) = \text{vec}(\mathbf{A})^T (\mathbf{D} \otimes \mathbf{B}) \text{vec}(\mathbf{C}). \quad (72)$$

Together with the update rule in Eq. (19), this concludes the proof of Theorem 7.  $\square$

## K Details on the simulation set-up

### K.1 Forward modeling

Populations of pyramidal neurons in the cortical gray matter are known to be the main drivers of the EEG signal [14, 19]. Here, we use a realistic volume conductor model of the human head to model the linear relationship between primary electrical source currents generated within these populations and the resulting scalp surface potentials captured by EEG electrodes. The lead field matrix,  $\mathbf{L} \in \mathbb{R}^{58 \times 2004}$ , which serves as the forward model in our simulations, was generated using the New York Head model [36]. The New York Head model provides a segmentation of an average human head into six different tissue types, taking into account the realistic anatomy and electrical tissue conductivities. In this model, 2004 dipolar current sources were placed evenly on the cortical

surface and 58 sensors were placed on the scalp according to the extended 10-20 system [94]. Finally, the lead field matrix was computed using the finite element method (FEM) for a given head geometry and exploiting the quasi-static approximation of Maxwell’s equations [14, 19, 36, 95].

Note that in accordance with the predominant orientation of pyramidal neuron assemblies, the orientation of all simulated source currents was fixed to be perpendicular to the cortical surface, so that only the scalar deflection of each source along the fixed orientation needs to be estimated. In real data analyses in Section 6 and Appendix L, however, surface normals are hard to estimate or even undefined in case of volumetric reconstructions. Consequently, we model each source in real data analyses as a full 3-dimensional current vector. This is achieved by introducing three variance parameters for each source within the source covariance matrix,  $\mathbf{\Gamma}^{3D} = \text{diag}(\gamma^{3D}) = [\gamma_1^x, \gamma_1^y, \gamma_1^z, \dots, \gamma_N^x, \gamma_N^y, \gamma_N^z]^\top$ . As all algorithms considered here model the source covariance matrix  $\mathbf{\Gamma}$  to be diagonal, this extension can be readily implemented. Correspondingly, a full 3D leadfield matrix,  $\mathbf{L}^{3D} \in \mathbb{R}^{M \times 3N}$ , is used.

## K.2 Pseudo-EEG signal generation

We simulated a sparse set of  $N_0 = 3$  active sources, which were placed at random positions on the cortex. To simulate the electrical neural activity of these sources,  $T = \{10, 20, 50, 100\}$  time points were sampled from a univariate linear autoregressive (AR) process, which models the activity at time  $t$  as a linear combination of the  $P$  past values:

$$x_i(t) = \sum_{p=1}^P a_i(p)x_i(t-p) + \xi_i(t), \text{ for } i = 1, 2, 3. \quad (73)$$

Here,  $a_i(p)$  for  $i = 1, 2, 3$  are linear AR coefficients, and  $P$  is the order of the AR model. The model residuals  $\xi_i(\cdot)$  for  $i = 1, 2, 3$  are also referred to as the innovation process; their variance determines the stability of the overall AR process. We here assume uncorrelated standard normal distributed innovations, which are independent for all sources. In this experiment, we use stable AR systems of order  $P = \{1, 2, 5, 7\}$ . The resulting source distribution, represented as  $\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(T)]$ , was projected to the EEG sensors through application of lead field matrix:  $\mathbf{Y}^{\text{signal}} = \mathbf{L}\mathbf{X}$ . Next we added Gaussian white noise to the sensor space signal. To this end, the same number of data points as the sources were sampled from a zero-mean normal distribution, where the time points assumed to be independent and identically distributed. The resulting noise distribution, represented as  $\mathbf{E} = [\mathbf{e}(1), \dots, \mathbf{e}(T)]$ , is then normalized by its Frobenius norm and added to the signal matrix  $\mathbf{Y}^{\text{signal}}$  as follows:  $\mathbf{Y} = \mathbf{Y}^{\text{signal}} + \frac{(1-\alpha)\|\mathbf{Y}^{\text{signal}}\|_F}{\alpha\|\mathbf{E}\|_F}\mathbf{E}$ , where  $\alpha$  determines the signal-to-noise ratio (SNR) in sensor space. Precisely, SNR is defined as follows:  $\text{SNR} = 20\log_{10}(\alpha/1-\alpha)$ . In this experiment the following values of  $\alpha$  were used:  $\alpha = \{0.55, 0.65, 0.7, 0.8\}$ , which correspond to the following SNRs:  $\text{SNR} = \{1.7, 5.4, 7.4, 12\}$  (dB). Interested reader can refer to Haufe and Ewald [37] for a more details on this simulation framework.

## K.3 Source reconstruction and evaluation metrics

We applied thin Dugh to the synthetic datasets described above. In addition to thin Dugh, one further Type-II Bayesian learning scheme, namely Champagne [33], and two Type-I source reconstruction schemes, namely S-FLEX [35] and eLORETA [34], were also included as benchmarks with respect to source reconstruction performance. S-FLEX is used as an example of a sparse Type-I Bayesian learning method based on  $\ell_1$ -norm minimization. As spatial basis functions, unit impulses were used, so that the resulting estimate was identical to the so-called minimum-current estimate (MCE) [96]. eLORETA estimate, as an example of a smooth inverse solution based on weighted  $\ell_2^2$ -norm minimization, was used with 5% regularization, whereas S-FLEX was fitted so that the residual variance was consistent with the ground-truth noise level. Note that the 5% rule is chosen as it gives the best performance across a subset of regularization values ranging between 0.5% to 15%. For thin Dugh, the noise variances as well as the variances of all voxels were initialized randomly by sampling from a standard normal distribution. The optimization program was terminated after reaching convergence. Convergence was defined if the relative change of the Frobenius-norm of the reconstructed sources between subsequent iterations was less than  $10^{-8}$ . A maximum of 1000 iterations was carried out if no convergence was reached beforehand.

Source reconstruction performance was evaluated according to two different measures, the *earth mover’s distance* (EMD), used to quantify the spatial localization accuracy, and the correlation

between the original and reconstructed sources,  $\hat{\mathbf{X}}$  and  $\mathbf{X}$ . The EMD metric measures the cost needed to map two probability distributions, defined on the same metric domain, into each other, see [21, 38]. It was applied here to the power of the true and estimated source activations defined on the cortical surface of the brain, which were obtained by taking the voxel-wise  $\ell_2$ -norm along the time domain. EMD was normalized to  $[0, 1]$ . The correlation between simulated and reconstructed source time courses was assessed as the mean of the absolute correlations obtained for each source, after optimally matching simulated and reconstructed sources. To this end, Pearson correlation between all pairs of simulated and reconstructed (i.e., those with non-zero activations) sources was measured. Each simulated source was matched to a reconstructed source based on maximum absolute correlation. Time-course correlation error (TCE) was then defined as one minus the average of these absolute correlations across sources. Each simulation was carried out 100 times using different instances of  $\mathbf{X}$  and  $\mathbf{E}$ , and the mean and standard error of the mean (SEM) of each performance measure across repetitions was calculated.

## L Real data analysis

### L.1 Auditory and visual evoked fields (AEF and VEF)

The MEG data used in this article were acquired in the Biomagnetic Imaging Laboratory at the University of California San Francisco (UCSF) with a CTF Omega 2000 whole-head MEG system from VSM MedTech (Coquitlam, BC, Canada) with 1200 Hz sampling rate. The neural responses for one subject’s auditory evoked fields (AEF) and visual evoked fields (VEF) were localized. The AEF response was elicited while subjects were passively listening to 600 ms duration tones (1 kHz) presented binaurally. Data from 120 trial epochs were analysed. The VEF response was elicited while subjects were viewing images of objects projected onto a screen and subjects were instructed to name the objects verbally. Both AEF and VEF data were first digitally filtered from 1 to 70 Hz to remove artifacts and DC offset, time-aligned to the stimulus. Different number of trials were included for algorithm analyses. The pre-stimulus window was selected to be  $-100$  ms to  $-5$  ms and the post-stimulus time window was selected to be 60 ms to 180 ms, where 0 ms is the onset of the tone. The lead field for each subject was calculated with NUTMEG [97] using a single-sphere head model (two spherical orientation lead fields) and an 8 mm voxel grid. Each column was normalized to have a norm of unity. Further details on these datasets can be found in [33, 98–100].

Figure 6 shows the reconstructed sources of the Auditory Evoked Fields (AEF) from a representative subject using Champagne, thin and full Dugh. In this case, we tested the reconstruction performance of all algorithms with the number of trials limited to 20 and 120. As Figure 6 demonstrates, the performance of Dugh remains robust as the number of trials is increased to 20 and 120 in Figure 6. Finally, the VEF performance of benchmark algorithm eLORETA is demonstrated in Figure 7.

### L.2 EEG Data: Faces vs scrambled pictures

A publicly available EEG dataset (128-channel Biosemi ActiveTwo system) was downloaded from the SPM website (<http://www.fil.ion.ucl.ac.uk/spm/data/mmfaces>) and the lead field was calculated in SPM8 using a three-shell spherical model at the coarse resolution of 5124 voxels at approximately 8 mm spacing. These EEG data were also obtained during a visual response paradigm that involved randomized presentation of at least 86 faces and 86 scrambled faces. To examine the differential responses to faces across all trials, the averaged responses to scrambled-faces were subtracted from the averaged responses to faces. The result is demonstrated in Figure 8.

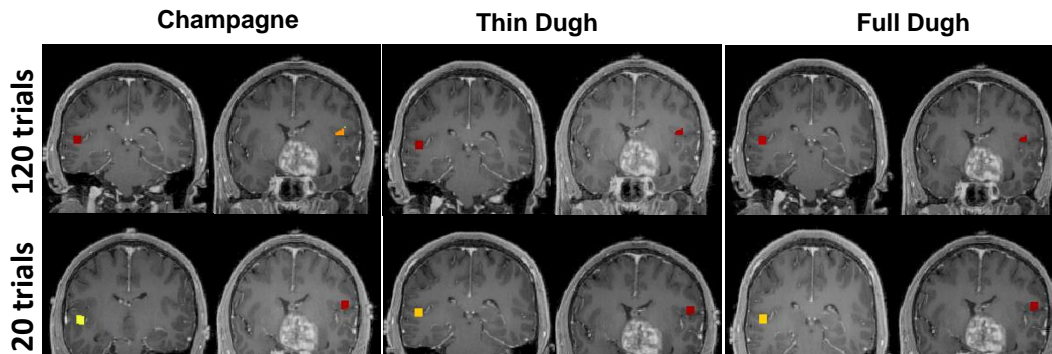


Figure 6: Robustness of Dugh and Champagne performance when the number of trials is increased to 20 and 120.

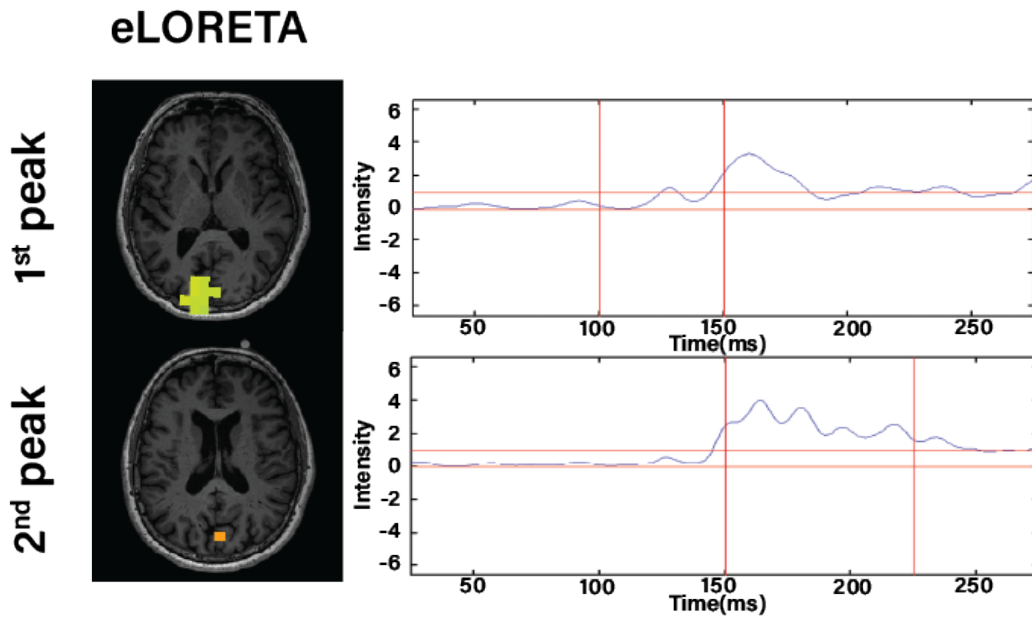


Figure 7: VEF performance of benchmark algorithm eLORETA. This benchmark did not yield reliable results for 5 trial epochs. Even when the number of trials were increased to 20, benchmark's performance yielded neither good spatial localization of the two visual cortical areas nor good estimation of the time courses of these activations.



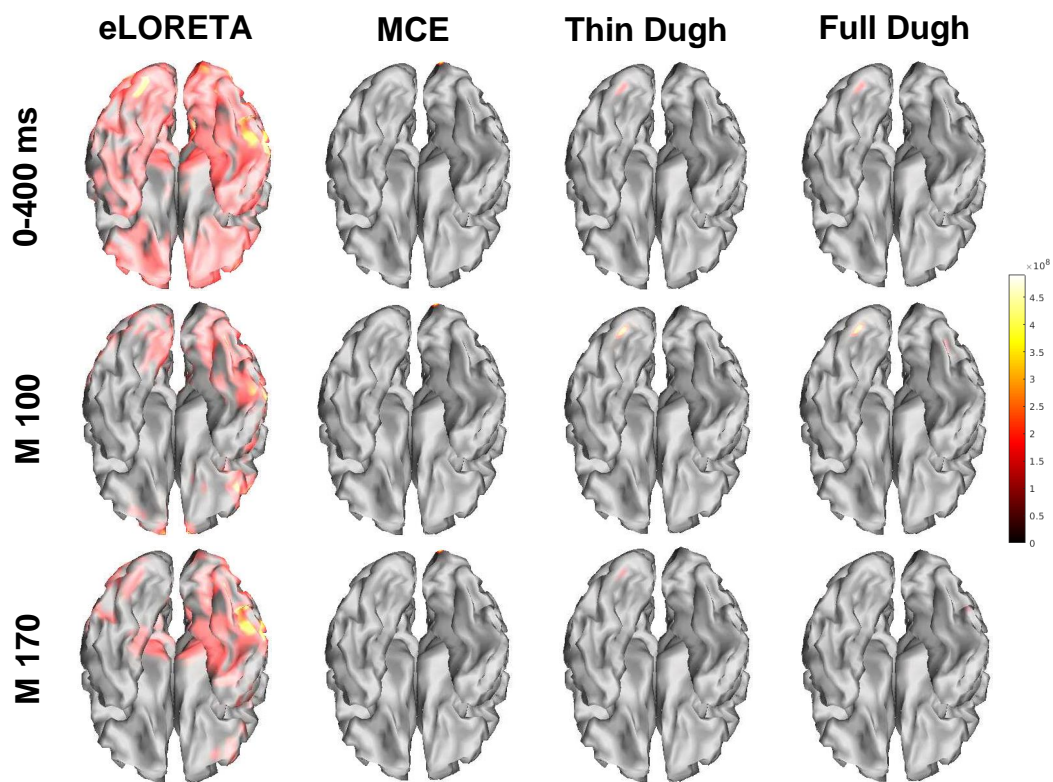


Figure 8: Performance of Dugh and benchmarks on EEG data acquired during a face recognition task. Dugh was able to provide more focal and distinct activations for the M100 and M170 responses that were not clearly identified using the benchmarks eLORETA and MCE.