

# Joint Learning of Full-structure Noise in Hierarchical Bayesian Regression Models

Ali Hashemi, Chang Cai, Yijing Gao, Sanjay Ghosh,

Klaus-Robert Müller\*, Member, IEEE, Srikantan S. Nagarajan\*, Fellow, IEEE, and Stefan Haufe\*

**Abstract**— We consider the reconstruction of brain activity from electroencephalography (EEG). This inverse problem can be formulated as a linear regression with independent Gaussian scale mixture priors for both the source and noise components. Crucial factors influencing accuracy of source estimation are not only the noise level but also its correlation structure, but existing approaches have not addressed estimation of noise covariance matrices with full structure. To address this shortcoming, we develop hierarchical Bayesian (type-II maximum likelihood) models for observations with latent variables for source and noise, which are estimated jointly from data. As an extension to classical sparse Bayesian learning (SBL), where across-sensor observations are assumed to be independent and identically distributed, we consider Gaussian noise with full covariance structure. Using the majorization-maximization framework and Riemannian geometry, we derive an efficient algorithm for updating the noise covariance along the manifold of positive definite matrices. We demonstrate that our algorithm has guaranteed and fast convergence and validate it in simulations and with real MEG data. Our results demonstrate that the novel framework significantly improves upon state-of-the-art techniques in the real-world scenario where the noise is indeed non-diagonal and fully-structured. Our method has applications in many domains beyond biomagnetic inverse problems.

**Index Terms**— EEG/MEG Brain Source Imaging, Hierarchical Bayesian Learning, Majorization Minimization, Sparse Bayesian Learning, Type-II Maximum-Likelihood.

## I. INTRODUCTION

HAVING precise knowledge of the noise distribution is a fundamental requirement for obtaining accurate solutions in many regression problems [1], particularly for biomedical imaging applications such as neural encoding models for task-based fMRI analyses [2]–[4] or magneto-

or electroencephalography (M/EEG) inverse problems [5]–[7]. In these biomedical imaging applications, however, it is impossible to separately estimate this noise distribution, as distinct “noise-only” (baseline) measurements are not feasible. An alternative, therefore, is to design estimators that jointly optimize over the regression coefficients as well as over parameters of the noise distribution. This has been pursued both in a (penalized) maximum-likelihood settings (here referred to as *Type-I* approaches) [5], [8], [9] as well as in hierarchical Bayesian settings (referred to as *Type-II*) [6], [7], [10], [11]. Most contributions in the literature are, however, limited to the estimation of only a scalar noise level (homoscedastic noise) or a diagonal noise covariance (i.e., independent between different measurements, heteroscedastic noise) [12]–[14]. Considering scalar or diagonal noise covariance is a limiting assumption in practice as the noise interference in many realistic scenarios are highly correlated across measurements; and thus, have non-trivial off-diagonal elements.

In this paper, we consider the problem of electromagnetic brain source imaging (BSI) as our main application. The goal of BSI is to reconstruct brain activity from magneto- or electroencephalography (M/EEG), which can be formulated as a sparse Bayesian learning (SBL) problem. Specifically, it can be cast as a linear Bayesian regression model with independent Gaussian scale mixture priors on the parameters and noise. Extending classical SBL approaches, we here consider Gaussian noise with full covariance structure. Prominent source of correlated noise in this context are, for example, eye blinks, heart beats, muscular artifacts and line noise. Other realistic examples for the need for such full-structure noise can be found in the areas of array processing [15] or direction of arrival (DOA) estimation [16]. Algorithms that can accurately estimate noise with full covariance structure are expected to achieve more accurate regression models and predictions in this setting. Therefore, our contribution in this paper consists in developing an efficient optimization algorithm for jointly estimating the posterior of regression parameters as well as the noise distribution. More specifically, we consider linear regression with Gaussian scale mixture priors on the parameters and a full-structure multivariate Gaussian noise. We cast the problem as a hierarchical Bayesian (Type-II maximum-likelihood) regression problem, in which the *source variance hyperparameters* and a *full-structural noise covariance matrix* are jointly estimated by maximizing the Bayesian evidence of the model. We derive an efficient algorithm for jointly

\*Corresponding authors.

Ali Hashemi and Stefan Haufe are with the Uncertainty, Inverse Modeling and Machine Learning Group, Technische Universität Berlin, Germany. Stefan Haufe is also with Physikalisch-Technische Bundesanstalt Braunschweig and Berlin, Germany, and Charité – Universitätsmedizin Berlin, Germany (Email: {hashemi,haufe}@tu-berlin.de).

Chang Cai, Yijing Gao, Sanjay Ghosh, and Srikantan S. Nagarajan are with Department of Radiology and Biomedical Imaging, University of California, San Francisco, CA, USA (Email: sri@ucsf.edu).

Klaus-Robert Müller is with Machine Learning Group, Technische Universität Berlin, Germany, BIFOLD – Berlin Institute for the Foundations of Learning and Data, Berlin, Germany, Department of Artificial Intelligence, Korea University, Seoul, South Korea, and also with Max Planck Institute for Informatics, Saarbrücken, Germany (Email: klaus-robert.mueller@tu-berlin.de).

estimating the source variances and noise covariance, along a Riemannian manifold of positive definite (P.D.) matrices.

The paper is organized as follows: In Section II, after reviewing the necessary background on Type-II Bayesian learning, we introduce our proposed algorithm. Simulation studies demonstrating significant improvement in source localization for EEG/MEG brain source imaging are presented in Section IV. Finally, Section VI concludes the paper.

## II. TYPE-II BAYESIAN REGRESSION

We consider the linear model  $\mathbf{Y} = \mathbf{L}\mathbf{X} + \mathbf{E}$ , in which a forward or design matrix,  $\mathbf{L} \in \mathbb{R}^{M \times N}$ , is mapped to the measurements,  $\mathbf{Y}$ , by a set of coefficients or source components,  $\mathbf{X}$ . Depending on the setting, the problem of estimating  $\mathbf{X}$  given  $\mathbf{L}$  and  $\mathbf{Y}$  is called an inverse problem in physics, a multi-task regression problem in machine learning, or a multiple measurement vector (MMV) recovery problem in signal processing [17]. Adopting a signal processing terminology, the *measurement matrix*  $\mathbf{Y} \in \mathbb{R}^{M \times T}$  captures the activity of  $M$  sensors at  $T$  time instants,  $\mathbf{y}(t) \in \mathbb{R}^{M \times 1}, t = 1, \dots, T$ , while the *source matrix*,  $\mathbf{X} \in \mathbb{R}^{N \times T}$ , consists of the unknown activity of  $N$  sources at the same time instants,  $\mathbf{x}(t) \in \mathbb{R}^{N \times 1}, t = 1, \dots, T$ . The matrix  $\mathbf{E} = [\mathbf{e}(1), \dots, \mathbf{e}(T)] \in \mathbb{R}^{M \times T}$  represents  $T$  time instances of zero-mean Gaussian noise with full covariance  $\Lambda$ ,  $\mathbf{e}(t) \in \mathbb{R}^{M \times 1} \sim \mathcal{N}(0, \Lambda), t = 1, \dots, T$ , which is assumed to be independent of the source activations.

In this paper, we focus on M/EEG based brain source imaging (BSI). However, the proposed algorithm can be used in general regression settings, in particular for sparse signal recovery [18], [19] with a wide range of applications [20]. The goal of BSI is to infer the underlying brain activity  $\mathbf{X}$  from the EEG/MEG measurement  $\mathbf{Y}$  given a known forward operator, called lead field matrix  $\mathbf{L}$ . In practice,  $\mathbf{L}$  can be computed using discretization methods such as the finite element method (FEM) for a given head geometry and known electrical conductivities using the quasi-static approximation of Maxwell's equations [21], [22]. As the number of sensors is typically much smaller than the number of locations of potential brain sources, this inverse problem is highly ill-posed. This problem is addressed by imposing prior distributions on the model parameters and adopting a Bayesian treatment. This can be performed either through Maximum-a-Posteriori (MAP) estimation (*Type-I Bayesian learning*) [23]–[27] or, when the model has unknown hyperparameters, through Type-II Maximum-Likelihood estimation (*Type-II Bayesian learning*) [28]–[32]. In this paper, we focus on Type-II Bayesian learning, which assumes a family of prior distributions  $p(\mathbf{X}|\Theta)$  parameterized by a set of hyperparameters  $\Theta$ . These hyperparameters can be learned from the data along with the model parameters using a hierarchical Bayesian approach [29], [33] through the maximum-likelihood principle:

$$\Theta^{\text{II}} := \arg \max_{\Theta} p(\mathbf{Y}|\Theta) = \arg \max_{\Theta} \int p(\mathbf{Y}|\mathbf{X}, \Theta)p(\mathbf{X}|\Theta)d\mathbf{X}.$$

Here we assume a zero-mean Gaussian prior with diagonal covariance  $\Gamma = \text{diag}(\gamma)$  for the underlying source distribution. That is,  $\mathbf{x}(t) \in \mathbb{R}^{N \times 1} \sim \mathcal{N}(0, \Gamma), t = 1, \dots, T$ ,

where  $\gamma = [\gamma_1, \dots, \gamma_N]^T$  contains  $N$  distinct unknown variances associated to  $N$  modeled brain sources. In the Type-II Bayesian learning framework, modeling independent sources through a diagonal covariance matrix leads to sparsity of the resulting source distributions, i.e., at the optimum, many of the estimated source variances are zero. This mechanism is known as *sparse Bayesian learning* (SBL) and is also closely related to the concept of *automatic relevance determination* (ARD) [29] and *kernel Fisher discriminant* (KFD) [28]. Just as most other approaches, SBL makes the simplifying assumption of statistical independence between time samples. This leads to the following expression for the distribution of the sources and measurements:

$$p(\mathbf{X}|\Gamma) = \prod_{t=1}^T p(\mathbf{x}(t)|\Gamma) = \prod_{t=1}^T \mathcal{N}(0, \Gamma) \quad (1)$$

$$p(\mathbf{Y}|\mathbf{X}) = \prod_{t=1}^T p(\mathbf{y}(t)|\mathbf{x}(t)) = \prod_{t=1}^T \mathcal{N}(\mathbf{L}\mathbf{x}(t), \Lambda). \quad (2)$$

The parameters of the Type-II model,  $\Theta$ , are the unknown source variances and the noise covariance, i.e.,  $\Theta = \{\Gamma, \Lambda\}$ . The unknown parameters  $\Gamma$  and  $\Lambda$  are optimized based on the current estimates of the source variances and noise covariance in an alternating iterative process. Given initial estimates of  $\Gamma$  and  $\Lambda$ , the posterior distribution of the sources is a Gaussian of the form [34]

$$p(\mathbf{X}|\mathbf{Y}, \Gamma, \Lambda) = \prod_{t=1}^T \mathcal{N}(\bar{\mathbf{x}}(t), \Sigma_x), \text{ where} \quad (3)$$

$$\bar{\mathbf{x}}(t) = \Gamma \mathbf{L}^\top (\Sigma_y)^{-1} \mathbf{y}(t) \quad (4)$$

$$\Sigma_x = \Gamma - \Gamma \mathbf{L}^\top (\Sigma_y)^{-1} \mathbf{L} \Gamma \quad (5)$$

$$\Sigma_y = \mathbf{L} \Gamma \mathbf{L}^\top + \Lambda. \quad (6)$$

The estimated posterior parameters  $\bar{\mathbf{x}}(t)$  and  $\Sigma_x$  are then in turn used to update  $\Gamma$  and  $\Lambda$  as the minimizers of the negative log of the marginal likelihood  $p(\mathbf{Y}|\Gamma, \Lambda)$ ,  $-\log p(\mathbf{Y}|\Gamma, \Lambda)$ , which is given by [35]:

$$\begin{aligned} \mathcal{L}^{\text{II}}(\Gamma, \Lambda) &= \log |\Sigma_y| + \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \Sigma_y^{-1} \mathbf{y}(t) \\ &= \log |\Lambda + \mathbf{L} \Gamma \mathbf{L}^\top| + \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top (\mathbf{L} \Gamma \mathbf{L}^\top + \Lambda)^{-1} \mathbf{y}(t), \end{aligned} \quad (7)$$

where  $|\cdot|$  denotes the determinant of a matrix. This process is repeated until convergence. Given the final solution of the hyperparameters  $\Theta^{\text{II}} = \{\Gamma^{\text{II}}, \Lambda^{\text{II}}\}$ , the posterior source distribution is obtained by plugging these estimates into (2)–(5).

## III. PROPOSED METHOD: FULL-STRUCTURE NOISE (FUN) LEARNING

Here we propose a novel and efficient algorithm, full-structure noise (FUN) learning, which is able to learn the full covariance structure of the noise jointly within the Bayesian Type-II regression framework. We adopt the SBL assumption for the sources, leading to  $\Gamma$ -updates previously described

in the BSI literature under the name Champagne [30]. As a novelty and main focus of this paper, we here equip the SBL framework with the capability to jointly learn full noise covariances by invoking efficient methods from Riemannian geometry, in particular the geometric mean.

Note that the Type-II cost function in (7) is non-convex and thus non-trivial to optimize. A number of iterative algorithms such as *majorization-minimization* (MM) [36] have been proposed to address this challenge. Following the MM scheme, we first construct convex surrogate functions that *majorizes*  $\mathcal{L}^{\text{II}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda})$  in each iteration of the optimization algorithm. Then, we show the minimization equivalence between the constructed majoring functions and (7). This result is presented in the following theorem:

**Theorem 1.** *Let  $\boldsymbol{\Lambda}^k$  and  $\boldsymbol{\Sigma}_y^k$  be fixed values obtained in the  $(k)$ -th iteration of the optimization algorithm minimizing  $\mathcal{L}^{\text{II}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda})$ . Then, optimizing the non-convex Type-II ML cost function in (7),  $\mathcal{L}^{\text{II}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda})$ , with respect to  $\boldsymbol{\Gamma}$  is equivalent to optimizing the following convex function, which majorizes (7):*

$$\mathcal{L}_{\text{source}}^{\text{conv}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda}^k) = \text{tr} \left( \mathbf{L}^\top (\boldsymbol{\Sigma}_y^k)^{-1} \mathbf{L} \boldsymbol{\Gamma} \right) + \text{tr}(\mathbf{M}_S^k \boldsymbol{\Gamma}^{-1}), \quad (8)$$

where  $\mathbf{M}_S^k$  is defined as:

$$\mathbf{M}_S^k := \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t) \bar{\mathbf{x}}^k(t)^\top. \quad (9)$$

Similarly, optimizing  $\mathcal{L}^{\text{II}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda})$  with respect to  $\boldsymbol{\Lambda}$  is equivalent to optimizing the following convex majorizing function:

$$\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda}) = \text{tr} \left[ (\boldsymbol{\Sigma}_y^k)^{-1} \boldsymbol{\Lambda} \right] + \text{tr}(\mathbf{M}_N^k \boldsymbol{\Lambda}^{-1}), \quad (10)$$

where  $\mathbf{M}_N^k$  is defined as:

$$\mathbf{M}_N^k := \frac{1}{T} \sum_{t=1}^T (\mathbf{y}(t) - \mathbf{L} \bar{\mathbf{x}}^k(t)) (\mathbf{y}(t) - \mathbf{L} \bar{\mathbf{x}}^k(t))^\top. \quad (11)$$

*Proof:* The proof is presented in Appendix D.

We continue by considering the optimization of the cost functions  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda})$  and  $\mathcal{L}_{\text{source}}^{\text{conv}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda}^k)$  with respect to  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\Gamma}$ , respectively. Note that in case of noise covariances with full structure, the solution of  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda})$  with respect to  $\boldsymbol{\Lambda}$  lies in the  $(M^2 - M)/2$  Riemannian manifold of P.D. matrices. This consideration enables us to invoke efficient methods from Riemannian geometry (see [37]), which ensures that the solution at each step of the optimization is contained within the lower-dimensional solution space. Specifically, in order to optimize for the noise covariance, the algorithm calculates the geometric mean between the previously obtained statistical model covariance,  $\boldsymbol{\Sigma}_y^k$ , and the empirical sensor-space residuals,  $\mathbf{M}_N^k$ , in each iteration. Regarding the solution of  $\mathcal{L}_{\text{source}}^{\text{conv}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda}^k)$ , note that we adopt the SBL assumption for the sources by imposing a diagonal structure on the source covariance matrix,  $\boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma})$ , where  $\boldsymbol{\gamma} = [\gamma_1, \dots, \gamma_N]^\top$ . The update rules obtained from this algorithm are presented in the following theorems:

**Theorem 2.** *The cost function  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda})$  is strictly geodesically convex with respect to the P.D. manifold, and*

*its minimum with respect to  $\boldsymbol{\Lambda}$  can be attained according to the following update rule:*

$$\boldsymbol{\Lambda}^{k+1} \leftarrow (\boldsymbol{\Sigma}_y^k)^{\frac{1}{2}} \left( (\boldsymbol{\Sigma}_y^k)^{-1/2} \mathbf{M}_N^k (\boldsymbol{\Sigma}_y^k)^{-1/2} \right)^{\frac{1}{2}} (\boldsymbol{\Sigma}_y^k)^{\frac{1}{2}}. \quad (12)$$

*Proof:* A detailed proof can be found in Appendix E. Moreover, a geometric representation of the geodesic path between the pair of matrices  $\{\boldsymbol{\Sigma}_y^k, \mathbf{M}_N^k\}$  on the P.D. manifold and the geometric mean between them, representing the update for  $\boldsymbol{\Lambda}^{k+1}$ , is provided in Fig. 7 in Appendix E.

**Theorem 3.** *Constraining  $\boldsymbol{\Gamma}$  in (8) to the set of diagonal matrices with nonnegative elements  $\mathcal{S}$ , i.e.,  $\mathcal{S} = \{\boldsymbol{\Gamma} \mid \boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma}) = \text{diag}([\gamma_1, \dots, \gamma_N]^\top), \gamma_n \geq 0, \text{ for } n = 1, \dots, N\}$ ,*

$$\boldsymbol{\Gamma}^{k+1} = \underset{\boldsymbol{\Gamma} \in \mathcal{S}, \boldsymbol{\Lambda} = \boldsymbol{\Lambda}^k}{\arg \min} \text{tr} \left( \mathbf{L}^\top (\boldsymbol{\Sigma}_y^k)^{-1} \mathbf{L} \boldsymbol{\Gamma} \right) + \text{tr}(\mathbf{M}_S^k \boldsymbol{\Gamma}^{-1}), \quad (13)$$

leads to the following update rule for the source variances:

$$\begin{aligned} \boldsymbol{\Gamma}^{k+1} &= \text{diag}(\boldsymbol{\gamma}^{k+1}), \text{ where,} \\ \gamma_n^{k+1} &\leftarrow \sqrt{\frac{[\mathbf{M}_S^k]_{n,n}}{[\mathbf{L}^\top (\boldsymbol{\Sigma}_y^k)^{-1} \mathbf{L}]_{n,n}}} = \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}^k(t))^2}{\mathbf{L}_{\cdot n}^\top (\boldsymbol{\Sigma}_y^k)^{-1} \mathbf{L}_{\cdot n}}} \\ &\text{for } n = 1, \dots, N, \end{aligned} \quad (14)$$

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

*Proof:* A detailed proof can be found in Appendix F. Convergence of the resulting algorithm is shown in the following theorem:

**Theorem 4.** *Optimizing the non-convex Type-II ML cost function in (7),  $\mathcal{L}^{\text{II}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda})$  with alternating update rules for  $\boldsymbol{\Lambda}$  and  $\boldsymbol{\Gamma}$  in (12) and (14) leads to an MM algorithm with convergence guarantees.*

*Proof:* A detailed proof can be found in Appendix G.

**Remark 1.** *Note that (14) is identical to the update rule of the Champagne algorithm [30]. Besides, various recent Type-II schemes for learning diagonal noise covariance matrices that are rooted in the concept of SBL [6], [7] can also be derived as special cases of FUN learning. Specifically, imposing diagonal structure on the noise covariance matrix for the FUN algorithm, i.e.,  $\boldsymbol{\Lambda} \in \mathcal{S}$ , results in the noise variance update rules derived in [7] for heteroscedastic, and in [6] for homoscedastic noise. Here, heteroscedasticity refers to the common situation that measurements are contaminated with non-uniform noise levels across channels, while homoscedasticity only accounts for uniform noise levels. We explicitly demonstrate the noted connection in Appendix H.*

**Remark 2.** *Note that, although FUN is limited to estimate a diagonal source covariance matrix, e.g.  $\boldsymbol{\Gamma} = \text{diag}(\boldsymbol{\gamma})$ , this assumption can be relaxed for special cases. One particular such setting is when the inverse of  $[\mathbf{L}^\top (\boldsymbol{\Sigma}_y^k)^{-1} \mathbf{L}]$  is well-defined. This condition is fulfilled whenever the rank of the lead field matrix  $\mathbf{L}$  is less than the number of sensors. In the context of BSI, this scenario, for example, occurs when a region-level lead field – instead of a voxel-level lead field – is*

used. Under this condition, an update rule similar to (12) can be obtained for the full-structural source covariance matrix:

$$\boldsymbol{\Gamma}^{k+1} \leftarrow (\mathbf{C}_S^k)^{\frac{1}{2}} \left( (\mathbf{C}_S^k)^{-1/2} \mathbf{M}_S^k (\mathbf{C}_S^k)^{-1/2} \right)^{\frac{1}{2}} (\mathbf{C}_S^k)^{\frac{1}{2}}, \quad (15)$$

where  $\mathbf{C}_S^k$  is defined as  $\mathbf{C}_S^k := \left( \mathbf{L}^\top (\boldsymbol{\Sigma}_y^k)^{-1} \mathbf{L} \right)^{-1}$ . For additional extensions to other scenarios, please see discussion.

**Remark 3.** The theoretical results presented in Section III have been obtained for the scalar setting, where the orientations of the dipolar brain source are assumed to be perpendicular to the surface of the cortex and, hence, only the scalar deflection of each source along the fixed orientation needs to be estimated. In real data, surface normals are hard to estimate or even undefined in case of volumetric reconstructions. Consequently, we model each source here as a full 3-dimensional current vector. This is achieved by introducing three variance parameters for each source within the source covariance matrix,  $\boldsymbol{\Gamma}^{3D} = [\gamma_1^x, \gamma_1^y, \gamma_1^z, \dots, \gamma_N^x, \gamma_N^y, \gamma_N^z]^\top$ . Further details is provided in Appendix C.

Summarizing, the FUN learning approach, just like Champagne and other SBL algorithms, assumes independent Gaussian sources with individual variances (thus, diagonal source covariances), which are updated through (14). Extending the classical SBL setting, which assumes the noise distribution to be known, FUN models noise with full covariance structure, which is updated using (12). Algorithm 1 in Appendix A summarizes the update rules used.

#### IV. NUMERICAL SIMULATIONS

In this section, we compare the performance of the proposed algorithm to variants employing simpler (hom- and heteroscedastic) noise models through an extensive set of simulations. Our simulation setting is an adoption of the EEG inverse problem, where brain activity is to be reconstructed from simulated pseudo-EEG data [38]. MATLAB code for producing the results in the simulation study is also released as an open source package in a publicly accessible GitHub repository: <https://github.com/AliHashemi-ai/FUN-Learning>.

##### A. Pseudo-EEG Signal Generation

**Forward Modeling:** Populations of pyramidal neurons in the cortical gray matter are known to be the main drivers of the EEG signal [21], [22]. 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,  $\mathbf{L} \in \mathbb{R}^{58 \times 2004}$ , was generated using the New York Head model [39] taking into account the realistic anatomy and electrical tissue conductivities of an average human head. In this model, 2004 dipolar current sources were placed evenly on the cortical surface and 58 sensors were considered. The lead field matrix,  $\mathbf{L} \in \mathbb{R}^{58 \times 2004}$  was computed using the finite element method [39]. Note that the orientation of all source currents was fixed to be perpendicular to the cortical surface, so that only scalar source amplitudes needed to be estimated.

**Source and Noise Model:** We simulated a sparse set of  $N_0 = 5$  active sources that were placed at random positions on the cortex. To simulate the electrical neural activity of these sources,  $T = 200$  identically and independently distributed (i.i.d) points were sampled from a Gaussian distribution, yielding sparse source activation vectors  $\mathbf{x}(t)$ . 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:  $\mathbf{Y}^{\text{signal}} = \mathbf{L}\mathbf{X}$ . Gaussian additive noise was randomly sampled from a zero-mean normal distribution with full covariance matrix  $\boldsymbol{\Lambda}$ :  $\mathbf{e}(t) \in \mathbb{R}^{M \times 1} \sim \mathcal{N}(0, \boldsymbol{\Lambda}), t = 1, \dots, T$ . This setting is further referred to as *full-structural noise*. Note that we also generated noise with diagonal covariance matrix, referred to as *heteroscedastic noise*, in order to investigate the effect of model violation on reconstruction performance. The noise matrix  $\mathbf{E} = [\mathbf{e}(1), \dots, \mathbf{e}(T)] \in \mathbb{R}^{M \times T}$  is normalized by its Frobenius norm and added to the signal matrix  $\mathbf{Y}^{\text{signal}}$  as follows:

$$\mathbf{Y} = \mathbf{Y}^{\text{signal}} + \frac{(1 - \alpha) \|\mathbf{Y}^{\text{signal}}\|_F}{\alpha \|\mathbf{E}\|_F} \mathbf{E}, \quad (16)$$

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 the subsequently described experiments the following values of  $\alpha$  were used:  $\alpha=\{0.3, 0.35, 0.4, 0.45, 0.5, 0.55, 0.65, 0.7, 0.8\}$ , which correspond to the following SNRs:  $\text{SNR}=\{-12, -7.4, -5.4, -3.5, -1.7, 0, 1.7, 3.5, 5.4, 7.4, 12\}$  (dB).

**Evaluation Metrics and Simulation Setup:** We applied the full-structural noise learning approach on the synthetic datasets described above to recover the locations and time courses of the active brain sources. In addition to our proposed approach, two further Type-II Bayesian learning schemes, namely homoscedastic and heteroscedastic Champagne [6], [7], were also included as benchmarks with respect to source reconstruction performance and noise covariance estimation accuracy.

Source reconstruction performance was evaluated according to the following metrics. First, the *earth mover's distance* (EMD) [25], [40] was used to quantify the spatial localization accuracy. The EMD measures the cost needed to transform two probability distributions defined on the same metric domain (in this case, distributions of the true and estimated sources defined in 3D Euclidean brain space) into each other. EMD scores were normalized to  $[0, 1]$ . Second, the error in the reconstruction of the source time courses was measured. To this end, Pearson correlation between all pairs of simulated and reconstructed (i.e., those with non-zero activations) source time courses was assessed as the mean of the absolute correlations obtained for each source, after optimally matching simulated and reconstructed sources based on maximal absolute correlation. We also report another metric for evaluating the localization error as the average Euclidean distance (EUCL) (in mm) between each simulated source and the best (in terms of absolute correlations) matching reconstructed source. For assessing the recovery of the true support, we also compute *F<sub>1</sub>-measure* scores [41], [42]:  $F_1 = 2 \times TP / P + TP + FP$ , where  $P$  denotes the number of true active sources, while  $TP$  and  $FP$

are the numbers of true and false positive predictions. Note that perfect support recovery, i.e.,  $F_1 = 1$ , is only achieved when there is a perfect correspondence between ground-truth and estimated support.

To evaluate the accuracy of the noise covariance matrix estimation, the following two metrics were calculated: the Pearson correlation measuring the structural similarity between original and reconstructed noise covariance matrices,  $\Lambda$  and  $\hat{\Lambda}$ , denoted by  $\Lambda^{\text{sim}}$ , and the normalized mean squared error (NMSE) between  $\Lambda$  and  $\hat{\Lambda}$ , defined as  $\text{NMSE} = \|\hat{\Lambda} - \Lambda\|_F^2 / \|\Lambda\|_F^2$ . Note that NMSE measures the reconstruction of the true scale of the noise covariance matrix, while  $\Lambda^{\text{sim}}$  is scale-invariant and hence only quantifies the overall structural similarity between simulated and estimated noise covariance matrices.

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. Convergence of the optimization programs for each run 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.

## B. Results

Fig. 1 shows two simulated datasets with five active sources in presence of full-structure noise (upper panel) as well as heteroscedastic noise (lower panel) at 0 dB SNR. Topographic maps depict the locations of the ground-truth active brain sources (first column) along with the source reconstruction result of three noise learning schemes assuming noise with homoscedastic (second column), heteroscedastic (third column), and full (fourth column) structure. For each algorithm, the estimated noise covariance matrix is also plotted above the topographic map. Source reconstruction performance was measured in terms of EMD and time course correlation (Corr), and is summarized in the table next to each panel. Besides, the accuracy of the noise covariance matrix reconstruction was measured on terms of  $\Lambda^{\text{sim}}$  and NMSE. Results are included in the same table.

Fig. 1 (upper panel) allows for a direct comparison of the estimated noise covariance matrices obtained from the three different noise learning schemes. It can be seen that FUN learning can better capture the overall structure of ground truth full-structure noise as evidenced by lower NMSE and similarity errors compared to the heteroscedastic and homoscedastic algorithm variants that are only able to recover a diagonal matrix while enforcing the off-diagonal elements to zero. This behaviour results in higher spatial and temporal accuracy (lower EMD and time course error) for FUN learning compared to competing algorithms assuming diagonal noise covariance. This advantage is also visible in the topographic maps.

The lower-panel of Fig. 1 presents analogous results for the setting where the noise covariance is generated according to a heteroscedastic model. Note that the superior spatial and temporal reconstruction performance of the heteroscedastic

noise learning algorithm compared to the full-structure scheme is expected here because the simulated ground truth noise is indeed heteroscedastic. The full-structure noise learning approach, however, provides fairly reasonable performance in terms of EMD, time course correlation (corr), and  $\Lambda^{\text{sim}}$ , although it is designed to estimate a full-structure noise covariance matrix. The convergence behaviour of all three noise learning variants is also illustrated in Fig. 1. Note that the full-structure noise learning approach eventually reaches lower negative log-likelihood values in both scenarios, namely full-structure and heteroscedastic noise.

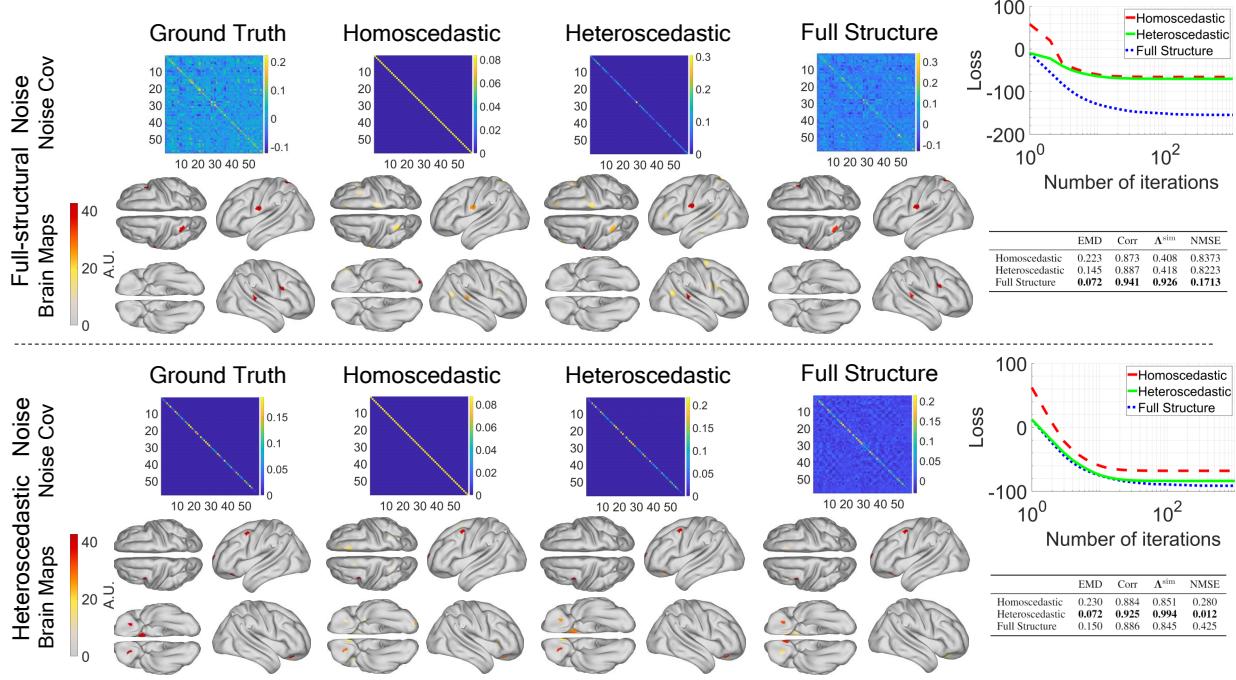
Fig. 2 shows the EMD, the time course reconstruction error, the EUCL and the F1 measure score incurred by three different noise learning approaches assuming homoscedastic (red), heteroscedastic (green) and full-structure (blue) noise covariances for a range of SNR values. The upper panel represents the evaluation metrics for the setting where the noise covariance is full-structure model, while the lower-panel depicts the same metric for simulated noise with heteroscedastic diagonal covariance. Concerning the first setting, FUN learning consistently outperforms its homoscedastic and heteroscedastic counterparts according to all evaluation metrics in particular in low-SNR settings. Consequently, as the SNR decreases, the gap between FUN learning and the two other variants increases. Conversely, heteroscedastic noise learning shows an improvement over FUN learning according to all evaluation metrics when the simulated noise is indeed heteroscedastic. However, note that the magnitude of this improvement is not as large as observed for the setting where the noise covariance is generated according to a full-structure model and then is estimated using the FUN approach.

Fig. 3 depicts the accuracy if the estimated noise covariance matrix reconstructed by three different noise learning approaches assuming noise with homoscedastic (red), heteroscedastic (green) and full (blue) structure. The ground truth noise covariance matrix either had full (upper row) or heteroscedastic (lower row) structure. Performance was measured in terms of similarity and NMSE. To be consistent with NMSE, we report “similarity error”, defined as  $1 - \Lambda^{\text{sim}}$ , instead of similarity,  $\Lambda^{\text{sim}}$ . Similar to the trend observed in Fig. 2, full-structure noise learning leads to better noise covariance estimation accuracy (lower NMSE and similarity error) for the full-structure noise model, while superior reconstruction performance is achieved for heteroscedastic noise learning when true noise covariance is heteroscedastic.

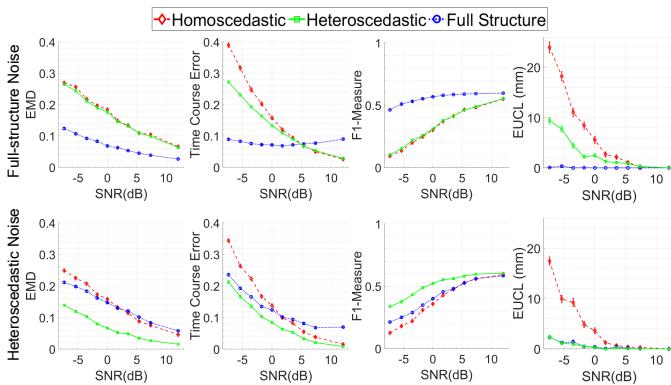
## V. ANALYSIS OF REAL MEG DATA

### A. Auditory and Visual Evoked Fields

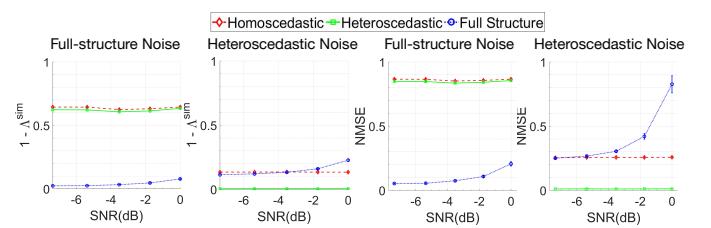
All MEG data used here were acquired in the Biomagnetic Imaging Laboratory at the University of California San Francisco (UCSF) with an Omega 2000 whole-head MEG system from CTF Inc. (Coquitlam, BC, Canada) at a sampling rate of 1200 Hz. All human participants provided informed written consent prior to study participation and received monetary compensation for their participation. The studies were approved by the University of California, San Francisco Committee on Human Research.



**Fig. 1:** Two examples of the simulated data with five active sources in presence of full-structural noise (upper panel) as well as heteroscedastic noise (lower panel) at 0 dB SNR. Topographic maps depict the locations of the ground-truth active brain sources (first column) along with the source reconstruction results of three noise learning schemes assuming noise with homoscedastic (second column), heteroscedastic (third column), or full structure (fourth column). For each algorithm, the estimated noise covariance matrix is also plotted above the topographic maps. The source reconstruction performance of these examples in terms of EMD and time course correlation (Corr) is summarized in the associated table next to each panel. Besides these two source reconstruction metrics, we also report the accuracy with which the ground-truth noise covariance was estimated in terms of the  $\Lambda^{\text{sim}}$  and NMSE metrics introduced above. The convergence behaviour of all three noise estimation approaches is also shown. Note that the full-structural noise learning approach converges to better minima of the negative log-likelihood than competing approaches regardless of whether the ground-truth noise covariance has full or heteroscedastic structure. However, an advantage in terms of reconstruction is only observed in the former case.



**Fig. 2:** Source reconstruction performance (mean  $\pm$  SEM) of the three different noise learning schemes for data generated by a realistic lead field matrix. Generated sensor signals were superimposed by either full-structure or heteroscedastic noise covering a wide range of SNRs. Performance was measured in terms of the earth mover's distance (EMD), time-course correlation error, F1-measure and Euclidean distance (EUCL) in (mm) between each simulated source and the reconstructed source with highest maximum absolute correlation.



**Fig. 3:** Accuracy of the noise covariance matrix reconstruction incurred by three different noise learning approaches assuming homoscedastic (red), heteroscedastic (green) and full-structural (blue) noise covariances. The ground-truth noise covariance matrix is either full-structure (upper row) or heteroscedastic diagonal (lower row). Performance is assessed in terms of the Pearson correlation between the entries of the original and reconstructed noise covariance matrices,  $\Lambda$  and  $\hat{\Lambda}$ , denoted by  $1 - \Lambda^{\text{sim}}$  (left column). Shown is the similarity error  $1 - \Lambda^{\text{sim}}$ . Further, the normalized mean squared error (NMSE) between  $\Lambda$  and  $\hat{\Lambda}$ , defined as  $\text{NMSE} = \|\hat{\Lambda} - \Lambda\|_F^2 / \|\Lambda\|_F^2$  is reported (right column).

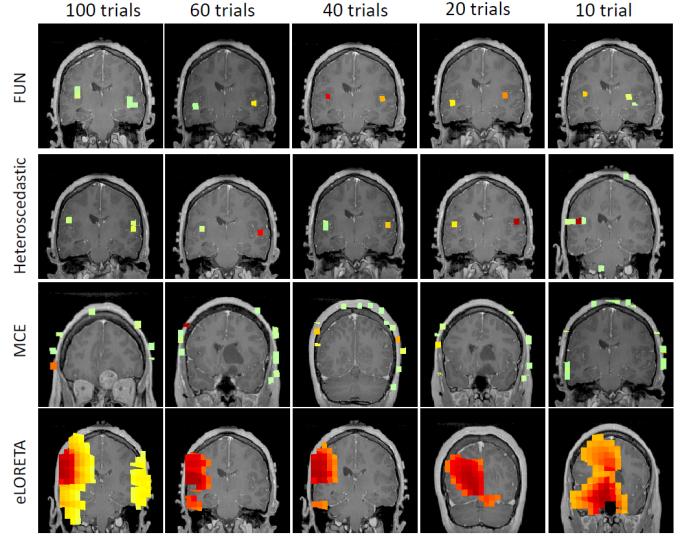
Lead-fields for each subject were calculated using NUT-MEG [43] assuming a single spherical shell volume conductor model resulting in only two spherical orientations ( $d_c = 2$ ). Lead-fields were constructed at a voxel resolution of 8 mm. Each lead-field column was normalized to have a norm of unity. Neural responses to auditory Evoked Fields (AEF) and visual evoked fields (VEF) stimulus were localized using the FUN algorithm and other benchmarks. The AEF response was elicited during passive listening to binaural tones (600 ms duration, carrier frequency of 1 kHz, 40 dB SL). The VEF response was elicited while subjects were viewing pictures of objects projected onto a screen and subjects were instructed to overtly name the objects [44], [45]. Up to 120 AEF and 100 VEF trials were collected. For both AEF and VEF data, trials with clear artifacts or visible noise in the MEG sensors that exceeded 10 pT fluctuations were excluded prior to source localization analysis.

Both AEF and VEF data were digitally filtered to a passband of 1 to 70 Hz and time-aligned to the stimulus onset. Averaging was then performed across sets of trials of increasing size: {10, 20, 40, 60, 100} trials for AEF, and {10, 20, 40} trials for VEF analyses. The pre-stimulus window was selected to be 100 ms prior to stimulus onset. The post-stimulus time window for AEF was selected to be +50 ms to +150 ms. For VEF data, we focused on source reconstruction in two time-windows – an early window ranging from +100 ms to +150 ms around the traditional M100 response, and a later time window ranging from +150 ms to +225 ms around the traditional M170 responses [35], [46]–[48].

Fig. 4 shows the reconstruction of the Auditory Evoked Fields (AEF) for different number of trial averages for a representative subject using FUN learning along with different Type-I and Type-II BSI benchmark methods. In addition to heteroscedastic Champagne, two classical non-SBL source reconstruction schemes were included for comparison. As an example of a sparse Type-I method based on  $\ell_1$ -norm minimization, the minimum-current estimate (MCE) algorithm [49] is shown. Additionally, eLORETA [50], representing a smooth inverse solution based on  $\ell_2$ -norm minimization, is also shown.

Reconstruction performance of all algorithms for different trial averaging with just 10, 20, 40, 60, and 100 trials are shown. All trials were selected randomly prior to averaging. As the subplots for different numbers of trial averages demonstrate, FUN learning can accurately localize bilateral auditory activity to Heschel's gyrus, the characteristic location of the primary auditory cortex, even with as few as 10 trials. In this challenging setting, all competing methods show inferior performance. These results highlight the importance of accurate noise covariance estimation on the fidelity of source reconstructions.

Fig. 5 shows the localization and time series reconstruction of visual evoked field (VEF) activity for a single subject using FUN and heteroscedastic noise learning Champagne, eLORETA and MCE. Reconstruction performance is again shown for the number of trials used for averaging ranging from 10, 20, and 40. Trials were randomly chosen from the full dataset without replacement prior to averaging. Within

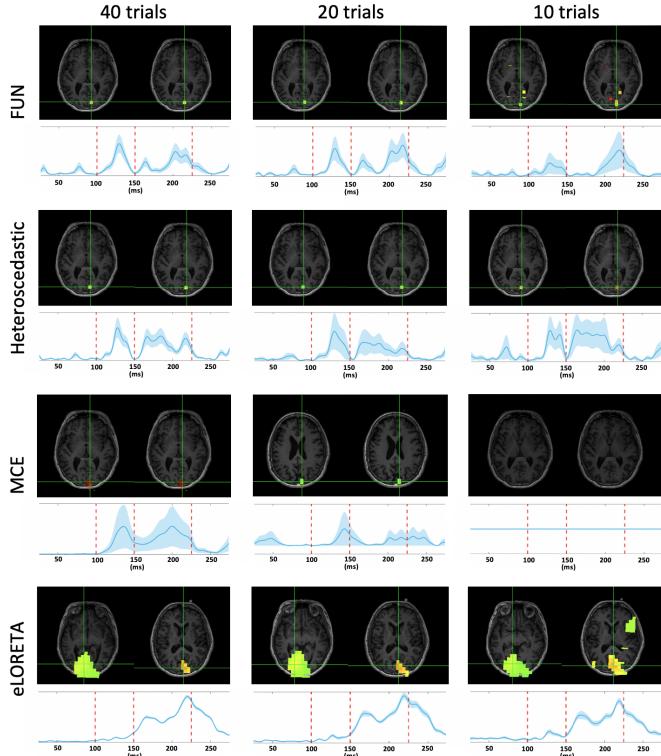


**Fig. 4:** Auditory evoked field (AEF) localization results from one representative subject for different numbers of trial averages using FUN learning, heteroscedastic Champagne, MCE and eLORETA. All reconstructions of FUN learning algorithm show focal sources at the expected locations of the auditory cortex. Even when limiting the number of trials to as few as 10 reconstruction result of FUN learning are accurate, while it severely affects the reconstruction performance of competing benchmarks methods.

each panel, the top shows the source localization of the M100 (1<sup>st</sup> peak) and M170 (2<sup>nd</sup> peak) responses, respectively. The time course of the source peak (indicated by the intersecting green lines) across a +25 ms to +275 ms window is presented below the source localization results. Blue lines represent the voxel power with arbitrary units averaged across ten independent experiments (that is, ten random selections of trials for trial averaging). Blue shades represent the standard error of the mean (SEM) across different trial averaging experiments. We also included three additional benchmark algorithms, sLORETA [51], S-FLEX [25] and the LCMV beamformer [52] in the supplementary material. In comparison with MCE and eLORETA benchmarks, FUN shows accurate localization capability, while both benchmarks did not yield reliable results for trial averaging as few as ten trials. Even when the number of trials used for averaging was increased to 20, the benchmarks yielded neither good spatial localization of the two visual cortical peaks, nor did they provide a reasonable estimation of the time courses of these activations. Furthermore, FUN detects two salient and clear peaks in each time-window in contrast to other benchmarks where the salience of the early and late peaks are less prominent. Results obtained from FUN are also robust across different SNRs/numbers of trial averages. For more benchmark results, please see supplementary Fig. 8.

### B. Resting-state data

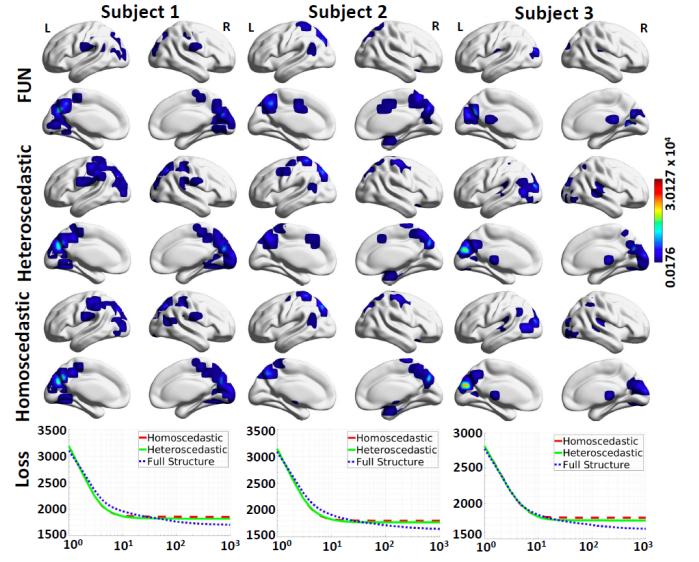
Resting-state data are particularly suited for the FUN algorithm because of the lack of baseline data on which the noise



**Fig. 5:** Localization and time series results of visual evoked field (VEF) activity for a single subject using FUN and benchmarks. Comparing with MCE and eLORETA, FUN shows accurate localization capability. Furthermore, FUN detects sharper 2<sup>nd</sup> peaks when compared to the heteroscedastic noise-learning Champagne, which is consistent with the sharp response of the VEF. The results obtained by FUN are robust across different SNRs/numbers of trial averages. For additional benchmark results, please see supplementary Fig. 8.

distribution could be estimated. Here, we show that FUN is able to learn the underlying noise distribution and consistently recover brain activity. For this analysis, three subjects were instructed simply to keep their eyes closed and remain awake. We collected four trials per subject, where each trial was one minute long. We randomly chose 30 seconds or equivalently 36000 time samples for brain source reconstruction from one trial of each subject. These resting-state MEG data were digitally filtered using a pass-band ranging from 8 to 12 Hz (alpha band) to remove artifacts and DC offset.

Localization of resting state alpha band activity from the three subjects are shown in Fig. 6. The first three columns show the estimated source covariance patterns (with the application of a threshold of 10% the peak value) for the three noise learning variants of Champagne. Each row represents one subject. The corresponding loss function values across 1000 iterations are shown in the last column. FUN consistently localizes all subjects' brain activity predominantly near the midline occipital lobe or posterior cingulate gyrus consistent with expected locations of alpha generators known to dominate resting-state activity.



**Fig. 6:** Localization of resting-state brain activity for three subjects using FUN and the heteroscedastic and homoscedastic noise learning variants of Champagne. The source variance patterns estimated by each algorithm are projected onto the cortical surface. The convergence behaviour of all three noise estimation approaches is also shown in terms of the negative log-likelihood cost function. FUN converges to better minima when compared to these benchmarks.

## VI. DISCUSSION

In this paper, we focused on sparse regression within the hierarchical Bayesian regression framework and its application in EEG/MEG brain source imaging. We proposed an efficient optimization algorithm for jointly estimating Gaussian regression parameter distributions as well as Gaussian noise distributions with full covariance structure within a hierarchical Bayesian framework. Using the Riemannian geometry of positive definite matrices, we derived an efficient algorithm for jointly estimating brain source variances and noise covariance. The benefits of our proposed framework were evaluated within an extensive set of experiments in the context of the electromagnetic brain source imaging inverse problem and showed significant improvement upon state-of-the-art techniques in the realistic scenario where the noise has full covariance structure. The practical performance of our method is further assessed through analyses of real auditory evoked fields (AEF), visual evoked fields (VEF) and resting-state MEG data.

In the context of BSI, [53] proposed a method for selecting a single regularization parameter based on cross-validation and maximum-likelihood estimation, while [54]–[58] assume more complex spatio-temporal noise covariance structures. A common limitation of these works is, however, that the noise level is not estimated as part of the source reconstruction problem on task-related data but from separate noise recordings. Our proposed algorithm substantially differs in this respect, as it learns the noise covariance jointly with the brain source distribution from the same data. This joint estimation perspective is opposed to a step-wise independent estimation process that can

cause to error accumulation. The idea of joint estimation of brain source activity and noise covariance has been previously proposed for Type-I learning methods in [5], [9]. [5] proposed a method to extend the group Lasso class of algorithms to multi-task learning, where the noise covariance is estimated using an eigenvalue fit to the empirical sensor space residuals defined in (11). In contrast, FUN learning uses Riemannian geometry principles, e.g., the geometric mean between the sensor space residuals defined in (11) and the previously obtained statistical model covariance,  $\Sigma_y^k$ . This enables us to robustly estimate the noise covariance as part of the model, in contrast to the method proposed in [5], which estimates the noise covariance solely based on the eigenvalues of the observed sensor space residuals. Furthermore, in contrast to these Type-I likelihood estimation methods, FUN is a Type-II method, which learns the prior source distribution as part of the model fitting. Type-II methods have been reported to yield results that are consistently superior to those of Type-I methods [6], [7], [47], [48]. Our numerical results show that the same holds also for FUN learning, which performs on par or better than existing variants from the Type-II family (including conventional Champagne) in this study.

Noise learning has also attracted attention in functional magnetic resonance imaging (fMRI) [2]–[4], where various models like matrix-normal (MN), factor analysis (FA), and Gaussian-process (GP) regression have been proposed. The majority of the noise learning algorithms in the fMRI literature rely on the EM framework, which is quite slow in practice [6] and has convergence guarantees only under certain restrictive conditions [36], [59]–[61]. In contrast to these existing approaches, our proposed framework not only applies to the models considered in these papers, but also benefits from theoretically proven convergence guarantees. To be more specific, we showed in this paper that FUN learning is an instance of the wider class of majorization-minimization (MM) framework, for which provable fast convergence is guaranteed. It is worth emphasizing our contribution within the MM optimization context as well. Unlike many other MM implementations, where surrogate functions are minimized using an iterative approach, our proposed algorithm is more efficient because it obtains a closed-form solution for the surrogate function in each step.

While being broadly applicable (please see Appendix B for a comprehensive list of potential applications), our approach is nevertheless also limited by a number of factors. Although Gaussian noise distributions are commonly justified, it would be desirable to also include more robust (e.g., heavy-tailed) non-Gaussian noise distributions in our framework. Another limitation is that the superior performance of the full-structure noise learning technique comes at the expense of higher computational complexity compared to the variants assuming homoscedastic or heteroscedastic structure. Besides, signals in real-world scenarios often lie in a lower-dimensional space compared to the original high-dimensional ambient space due to the correlations that exist in the data. Therefore, imposing physiologically plausible constraints on the noise model, e.g., low-rank, Toeplitz or Kronecker structure [62], [63], not only provides side information that can be leveraged for the

reconstruction but also reduces the computational cost in two ways: a) by reducing the number of parameters and b) by taking advantage of efficient implementations using circular embeddings and the fast Fourier transform [64], [65]. Exploring efficient ways to incorporate these structural assumptions within a Riemannian framework is another direction of our future work.

## REFERENCES

- [1] B. Rakitsch, C. Lippert, K. Borgwardt, and O. Stegle, “It is all in the noise: Efficient multi-task gaussian process inference with structured residuals,” in *Proceedings of the 26th International Conference on Neural Information Processing Systems-Volume 1*, 2013, pp. 1466–1474.
- [2] M. Cai, N. W. Schuck, J. W. Pillow, and Y. Niv, “A Bayesian method for reducing bias in neural representational similarity analysis,” in *Advances in Neural Information Processing Systems*, 2016, pp. 4951–4959.
- [3] M. Shvartsman, N. Sundaram, M. Aoi, A. Charles, T. Willke, and J. Cohen, “Matrix-normal models for fMRI analysis,” in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2018, pp. 1914–1923.
- [4] M. B. Cai, M. Shvartsman, A. Wu, H. Zhang, and X. Zhu, “Incorporating structured assumptions with probabilistic graphical models in fMRI data analysis,” *Neuropsychologia*, p. 107500, 2020.
- [5] Q. Bertrand, M. Massias, A. Gramfort, and J. Salmon, “Handling correlated and repeated measurements with the smoothed multivariate square-root Lasso,” in *Advances in Neural Information Processing Systems*, 2019, pp. 3959–3970.
- [6] A. Hashemi, C. Cai, G. Kutyniok, K.-R. Müller, S. Nagarajan, and S. Haufe, “Unification of sparse Bayesian learning algorithms for electromagnetic brain imaging with the majorization minimization framework,” *NeuroImage*, vol. 239, p. 118309, 2021.
- [7] C. Cai, A. Hashemi, M. Diwakar, S. Haufe, K. Sekihara, and S. S. Nagarajan, “Robust estimation of noise for electromagnetic brain imaging with the Champagne algorithm,” *NeuroImage*, vol. 225, p. 117411, 2021.
- [8] H. B. Petersen and P. Jung, “Robust instance-optimal recovery of sparse signals at unknown noise levels,” *arXiv preprint arXiv:2008.08385*, 2020.
- [9] M. Massias, O. Fercoq, A. Gramfort, and J. Salmon, “Generalized concomitant multi-task lasso for sparse multimodal regression,” in *International Conference on Artificial Intelligence and Statistics*, 2018, pp. 998–1007.
- [10] D. P. Wipf and B. D. Rao, “An empirical Bayesian strategy for solving the simultaneous sparse approximation problem,” *IEEE Transactions on Signal Processing*, vol. 55, no. 7, pp. 3704–3716, 2007.
- [11] Z. Zhang and B. D. Rao, “Sparse signal recovery with temporally correlated source vectors using sparse Bayesian learning,” *IEEE Journal of Selected Topics in Signal Processing*, vol. 5, no. 5, pp. 912–926, 2011.
- [12] S. Van de Geer, J. Lederer *et al.*, “The Lasso, correlated design, and improved oracle inequalities,” in *From Probability to Statistics and Back: High-Dimensional Models and Processes—A Festschrift in Honor of Jon A. Wellner*. Institute of Mathematical Statistics, 2013, pp. 303–316.
- [13] A. Dalalyan, M. Hebiri, K. Meziani, and J. Salmon, “Learning heteroscedastic models by convex programming under group sparsity,” in *International Conference on Machine Learning*, 2013, pp. 379–387.
- [14] J. Lederer and C. L. Muller, “Don’t fall for tuning parameters: tuning-free variable selection in high dimensions with the TREX,” in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015, pp. 2729–2735.
- [15] T. Li and A. Nehorai, “Maximum likelihood direction finding in spatially colored noise fields using sparse sensor arrays,” *IEEE Transactions on Signal Processing*, vol. 59, no. 3, pp. 1048–1062, 2010.
- [16] C. E. Chen, F. Lorenzelli, R. E. Hudson, and K. Yao, “Stochastic maximum-likelihood DOA estimation in the presence of unknown nonuniform noise,” *IEEE Transactions on Signal Processing*, vol. 56, no. 7, pp. 3038–3044, 2008.
- [17] S. F. Cotter, B. D. Rao, K. Engan, and K. Kreutz-Delgado, “Sparse solutions to linear inverse problems with multiple measurement vectors,” *IEEE Transactions on Signal Processing*, vol. 53, no. 7, pp. 2477–2488, 2005.
- [18] E. J. Candès, J. K. Romberg, and T. Tao, “Stable signal recovery from incomplete and inaccurate measurements,” *Communications on pure and applied mathematics*, vol. 59, no. 8, pp. 1207–1223, 2006.
- [19] D. L. Donoho, “Compressed sensing,” *IEEE Transactions on Information Theory*, vol. 52, no. 4, pp. 1289–1306, 2006.

- [20] D. Malioutov, M. Cetin, and A. S. Willsky, "A sparse signal reconstruction perspective for source localization with sensor arrays," *IEEE Transactions on Signal Processing*, vol. 53, no. 8, pp. 3010–3022, 2005.
- [21] M. Hämäläinen, R. Hari, R. J. Ilmoniemi, J. Knuutila, and O. V. Lounasmaa, "Magnetoencephalography—theory, instrumentation, and applications to noninvasive studies of the working human brain," *Reviews of modern Physics*, vol. 65, no. 2, p. 413, 1993.
- [22] S. Baillet, J. C. Mosher, and R. M. Leahy, "Electromagnetic brain mapping," *IEEE Signal Processing Magazine*, vol. 18, no. 6, pp. 14–30, 2001.
- [23] R. D. Pascual-Marqui, C. M. Michel, and D. Lehmann, "Low resolution electromagnetic tomography: a new method for localizing electrical activity in the brain," *International Journal of psychophysiology*, vol. 18, no. 1, pp. 49–65, 1994.
- [24] I. F. Gorodnitsky, J. S. George, and B. D. Rao, "Neuromagnetic source imaging with FOCUSS: a recursive weighted minimum norm algorithm," *Electroencephalography and Clinical Neurophysiology*, vol. 95, no. 4, pp. 231–251, 1995.
- [25] S. Haufe, V. V. Nikulin, A. Ziehe, K.-R. Müller, and G. Nolte, "Combining sparsity and rotational invariance in EEG/MEG source reconstruction," *NeuroImage*, vol. 42, no. 2, pp. 726–738, 2008.
- [26] A. Gramfort, M. Kowalski, and M. Hämäläinen, "Mixed-norm estimates for the M/EEG inverse problem using accelerated gradient methods," *Physics in Medicine and Biology*, vol. 57, no. 7, p. 1937, 2012.
- [27] S. Castaño-Candamil, J. Höhne, J.-D. Martínez-Vargas, X.-W. An, G. Castellanos-Domínguez, and S. Haufe, "Solving the EEG inverse problem based on space–time–frequency structured sparsity constraints," *NeuroImage*, vol. 118, pp. 598–612, 2015.
- [28] S. Mika, G. Rätsch, and K.-R. Müller, "A mathematical programming approach to the kernel fisher algorithm," *Advances in Neural Information Processing Systems*, vol. 13, pp. 591–597, 2001.
- [29] M. E. Tipping, "Sparse Bayesian learning and the relevance vector machine," *Journal of Machine Learning Research*, vol. 1, no. Jun, pp. 211–244, 2001.
- [30] D. Wipf and S. Nagarajan, "A unified Bayesian framework for MEG/EEG source imaging," *NeuroImage*, vol. 44, no. 3, pp. 947–966, 2009.
- [31] M. W. Seeger and D. P. Wipf, "Variational Bayesian inference techniques," *IEEE Signal Processing Magazine*, vol. 27, no. 6, pp. 81–91, 2010.
- [32] W. Wu, S. Nagarajan, and Z. Chen, "Bayesian machine learning: EEG\MEG signal processing measurements," *IEEE Signal Processing Magazine*, vol. 33, no. 1, pp. 14–36, 2016.
- [33] D. P. Wipf and B. D. Rao, "Sparse Bayesian learning for basis selection," *IEEE Transactions on Signal Processing*, vol. 52, no. 8, pp. 2153–2164, 2004.
- [34] K. Sekihara and S. S. Nagarajan, *Electromagnetic brain imaging: a Bayesian perspective*. Springer, 2015.
- [35] D. P. Wipf, J. P. Owen, H. T. Attias, K. Sekihara, and S. S. Nagarajan, "Robust Bayesian estimation of the location, orientation, and time course of multiple correlated neural sources using MEG," *NeuroImage*, vol. 49, no. 1, pp. 641–655, 2010.
- [36] Y. Sun, P. Babu, and D. P. Palomar, "Majorization-minimization algorithms in signal processing, communications, and machine learning," *IEEE Transactions on Signal Processing*, vol. 65, no. 3, pp. 794–816, 2017.
- [37] P. Petersen, S. Axler, and K. Ribet, *Riemannian geometry*. Springer, 2006, vol. 171.
- [38] S. Haufe and A. Ewald, "A simulation framework for benchmarking EEG-based brain connectivity estimation methodologies," *Brain topography*, pp. 1–18, 2016.
- [39] Y. Huang, L. C. Parra, and S. Haufe, "The New York head — a precise standardized volume conductor model for EEG source localization and tES targeting," *NeuroImage*, vol. 140, pp. 150–162, 2016.
- [40] Y. Rubner, C. Tomasi, and L. J. Guibas, "The earth mover's distance as a metric for image retrieval," *International Journal of Computer Vision*, vol. 40, no. 2, pp. 99–121, 2000.
- [41] N. Chinchor and B. M. Sundheim, "Muc-5 evaluation metrics," in *Fifth Message Understanding Conference (MUC-5)*, 1993.
- [42] C. van Rijsbergen, "Information retrieval," 1979.
- [43] S. S. Dalal, J. Zumer, V. Agrawal, K. Hild, K. Sekihara, and S. Nagarajan, "NUTMEG: a neuromagnetic source reconstruction toolbox," *Neurology & Clinical Neurophysiology: NCN*, vol. 2004, p. 52, 2004.
- [44] L. B. Hinkley, C. L. Dale, T. L. Luks, A. M. Findlay, P. Bukshpun, N. Pojman, T. Thieu, W. K. Chung, J. Berman, T. P. Roberts *et al.*, "Sensorimotor cortical oscillations during movement preparation in 16p11. 2 deletion carriers," *Journal of Neuroscience*, vol. 39, no. 37, pp. 7321–7331, 2019.
- [45] L. B. Hinkley, E. De Witte, M. Cahill-Thompson, D. Mizuri, C. Garrett, S. Honma, A. Findlay, M. L. Gorno-Tempini, P. Tarapore, H. E. Kirsch *et al.*, "Optimizing magnetoencephalographic imaging estimation of language lateralization for simpler language tasks," *Frontiers in Human Neuroscience*, vol. 14, p. 105, 2020.
- [46] S. S. Dalal, J. M. Zumer, A. G. Guggisberg, M. Trumpis, D. D. Wong, K. Sekihara, and S. S. Nagarajan, "MEG/EEG source reconstruction, statistical evaluation, and visualization with NUTMEG," *Computational Intelligence and Neuroscience*, vol. 2011, 2011.
- [47] J. P. Owen, D. P. Wipf, H. T. Attias, K. Sekihara, and S. S. Nagarajan, "Performance evaluation of the Champagne source reconstruction algorithm on simulated and real M/EEG data," *Neuroimage*, vol. 60, no. 1, pp. 305–323, 2012.
- [48] C. Cai, M. Diwakar, D. Chen, K. Sekihara, and S. S. Nagarajan, "Robust empirical Bayesian reconstruction of distributed sources for electromagnetic brain imaging," *IEEE Transactions on Medical Imaging*, vol. 39, no. 3, pp. 567–577, 2019.
- [49] K. Matsuura and Y. Okabe, "Selective minimum-norm solution of the biomagnetic inverse problem," *IEEE Transactions on Biomedical Engineering*, vol. 42, no. 6, pp. 608–615, 1995.
- [50] R. D. Pascual-Marqui, "Discrete, 3D distributed, linear imaging methods of electric neuronal activity. Part 1: exact, zero error localization," 2007.
- [51] R. D. Pascual-Marqui *et al.*, "Standardized low-resolution brain electromagnetic tomography (sLORETA): technical details," *Methods Find Exp Clin Pharmacol*, vol. 24, no. Suppl D, pp. 5–12, 2002.
- [52] B. D. Van Veen, W. Van Drongelen, M. Yuchtman, and A. Suzuki, "Localization of brain electrical activity via linearly constrained minimum variance spatial filtering," *IEEE Transactions on Biomedical Engineering*, vol. 44, no. 9, pp. 867–880, 1997.
- [53] D. A. Engemann and A. Gramfort, "Automated model selection in covariance estimation and spatial whitening of MEG and EEG signals," *NeuroImage*, vol. 108, pp. 328–342, 2015.
- [54] H. M. Huijzen, J. C. De Munck, L. J. Waldorp, and R. P. Grasman, "Spatiotemporal EEG/MEG source analysis based on a parametric noise covariance model," *IEEE Transactions on Biomedical Engineering*, vol. 49, no. 6, pp. 533–539, 2002.
- [55] J. C. De Munck, H. M. Huijzen, L. J. Waldorp, and R. Heethaar, "Estimating stationary dipoles from MEG/EEG data contaminated with spatially and temporally correlated background noise," *IEEE Transactions on Signal Processing*, vol. 50, no. 7, pp. 1565–1572, 2002.
- [56] F. Bijma, J. C. De Munck, H. M. Huijzen, and R. M. Heethaar, "A mathematical approach to the temporal stationarity of background noise in MEG/EEG measurements," *NeuroImage*, vol. 20, no. 1, pp. 233–243, 2003.
- [57] J. C. De Munck, F. Bijma, P. Gaura, C. A. Sieluzycki, M. I. Branco, and R. M. Heethaar, "A maximum-likelihood estimator for trial-to-trial variations in noisy MEG/EEG data sets," *IEEE Transactions on Biomedical Engineering*, vol. 51, no. 12, pp. 2123–2128, 2004.
- [58] S. C. Jun, S. M. Plis, D. M. Ranken, and D. M. Schmidt, "Spatiotemporal noise covariance estimation from limited empirical magnetoencephalographic data," *Physics in Medicine & Biology*, vol. 51, no. 21, p. 5549, 2006.
- [59] T. T. Wu, K. Lange *et al.*, "The MM alternative to EM," *Statistical Science*, vol. 25, no. 4, pp. 492–505, 2010.
- [60] M. W. Jacobson and J. A. Fessler, "An expanded theoretical treatment of iteration-dependent majorize-minimize algorithms," *IEEE Transactions on Image Processing*, vol. 16, no. 10, pp. 2411–2422, 2007.
- [61] M. Razaviyayn, M. Hong, and Z.-Q. Luo, "A unified convergence analysis of block successive minimization methods for nonsmooth optimization," *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 1126–1153, 2013.
- [62] A. Breloy, Y. Sun, P. Babu, G. Ginolhac, and D. P. Palomar, "Robust rank constrained kronecker covariance matrix estimation," in *2016 50th Asilomar Conference on Signals, Systems and Computers*. IEEE, 2016, pp. 810–814.
- [63] B. Xin, Y. Wang, W. Gao, and D. Wipf, "Building invariances into sparse subspace clustering," *IEEE Transactions on Signal Processing*, vol. 66, no. 2, pp. 449–462, 2017.
- [64] P. Babu, "MELT—maximum-likelihood estimation of low-rank Toeplitz covariance matrix," *IEEE Signal Processing Letters*, vol. 23, no. 11, pp. 1587–1591, 2016.
- [65] A. Hashemi, Y. Gao, C. Cai, S. Ghosh, K. R. Müller, S. S. Nagarajan, and S. Haufe, "Efficient hierarchical Bayesian inference for spatio-temporal regression models in neuroimaging," in *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.

## Supplementary Material

### APPENDIX

#### A. Summary of the Full-structure noise (FUN) learning method

**Algorithm 1:** Full-structure noise (FUN) learning

**Input:** The lead field matrix  $\mathbf{L} \in \mathbb{R}^{M \times N}$  and the measurement vectors  $\mathbf{y}(t) \in \mathbb{R}^{M \times 1}, t = 1, \dots, T$ .  
**Result:** The estimated prior source variances  $[\gamma_1, \dots, \gamma_N]^\top$ , noise covariance  $\Lambda$ , the posterior mean  $\bar{\mathbf{x}}(t)$  and covariance  $\Sigma_{\mathbf{x}}$  of the sources.

- 1 Choose a random initial value for  $\Lambda$  as well as  $\gamma = [\gamma_1, \dots, \gamma_N]^\top$ , and construct  $\Gamma = \text{diag}(\gamma)$ .
- 2 Calculate the statistical covariance  $\Sigma_y = \Lambda + \mathbf{L}\Gamma\mathbf{L}^\top$ .
- 3 Initialize  $k \leftarrow 1$
- 4 **Repeat**
- 5     Calculate the posterior mean as  $\bar{\mathbf{x}}(t) = \Gamma\mathbf{L}^\top(\Sigma_y)^{-1}\mathbf{y}(t)$ .
- 6     Calculate  $\mathbf{M}_N^k$  based on (11), and update  $\Lambda$  based on (12).
- 7     Calculate  $\mathbf{M}_S^k$  based on (9), and update  $\Gamma$  and  $\gamma_n$  for  $n = 1, \dots, N$  based on (14).
- 8      $k \leftarrow k + 1$
- 9 **Until** stopping condition is satisfied:  $\|\bar{\mathbf{x}}^{k+1} - \bar{\mathbf{x}}^k\|_2^2 \leq \epsilon$  or  $k = k_{max}$ ;
- 10 Calculate the posterior covariance as  $\Sigma_{\mathbf{x}} = \Gamma - \Gamma\mathbf{L}^\top(\Sigma_y)^{-1}\mathbf{L}\Gamma$ .

#### B. Broader Impact

Although this paper focuses on electromagnetic brain source imaging, our proposed algorithm is suitable for a wider range of applications. The same concepts used here for full-structure noise learning could be employed in other contexts where hyperparameters like kernel widths in Gaussian process regression [66] or dictionary elements in the dictionary learning problem [67] need to be inferred from data. The FUN learning algorithm may also prove useful for practical scenarios in which model residuals are expected to be correlated, e.g., probabilistic canonical correlation analysis (CCA) [68], spectral independent component analysis (ICA) [69], direction of arrival (DoA) and channel estimation in massive Multiple Input Multiple Output (MIMO) systems [70]–[72], robust portfolio optimization in finance [73], covariance matching and estimation [74]–[80], graph learning [81], thermal field reconstruction [82]–[84], and brain functional imaging [85]. It is also straightforward to incorporate our optimization procedure within more complex models with hierarchical priors. This includes problems that optimize the evidence lower bound (ELBO) cost function in variational Bayesian inference [86] or variational autoencoders [87].

#### C. Recovering brain sources with free orientation

The theoretical results presented in Section III have been obtained for the scalar setting, where the orientations of the dipolar brain source are assumed to be perpendicular to the surface of the cortex and, hence, only the scalar deflection of each source along the fixed orientation needs to be estimated. In real data, surface normals are hard to estimate or even undefined in case of volumetric reconstructions. Consequently, we model each source here as a full 3-dimensional current vector. This is achieved by introducing three variance parameters for each source within the source covariance matrix,  $\Gamma^{3D} = \text{diag}(\gamma^{3D}) = [\gamma_1^x, \gamma_1^y, \gamma_1^z, \dots, \gamma_N^x, \gamma_N^y, \gamma_N^z]^\top$ . Correspondingly, a full 3D leadfield matrix,  $\mathbf{L}^{3D} \in \mathbb{R}^{M \times 3N}$ , is used. More specifically, we define  $\mathbf{L}^{3D} = [\mathbf{L}_1, \dots, \mathbf{L}_N]$ , where  $N$  is the number of voxels under consideration and  $\mathbf{L}_n = [\mathbf{L}_n^1, \dots, \mathbf{L}_n^{d_c}] \in \mathbb{R}^{M \times d_c}$  is the leadfield matrix for  $n$ -th voxel with  $d_c$  orientations. The  $k$ -th column of  $\mathbf{L}_n$ , i.e.  $\mathbf{L}_n^k$  for  $k = 1, \dots, d_c$ , represents the signal vector that would be observed at the scalp given a unit current source or dipole at the  $n$ -th voxel with a fixed orientation in the  $k$ -th direction. The voxel dimension  $d_c$  is commonly set to 3 for EEG, and MEG with realistic volume conductor models, and 2 for MEG with single spherical shell models.

In this scenario,  $\mathbf{x}_n(t) = [x_n^1(t), \dots, x_n^{d_c}(t)]^\top \in \mathbb{R}^{d_c \times 1}$  models the  $n$ -th voxel intensity at time  $t$ , which we assume it with  $d_c$  orientations. Then, the generative probabilistic model for the sensor data at time point  $t$  can be written as:

$$\mathbf{y}(t) = \mathbf{L}^{3D}\mathbf{x}(t) + \mathbf{e}(t) = \sum_{n=1}^N \mathbf{L}_n \mathbf{x}_n(t) + \mathbf{e}(t),$$

with prior distributions  $\mathbf{x}(t) \sim \mathcal{N}(\mathbf{0}, \Gamma^{3D})$ , where  $\Gamma^{3D}$  is defined as  $d_c N \times d_c N$  block diagonal matrix expressed as

$$\Gamma^{3D} = \begin{bmatrix} \gamma_1 \mathbf{I}_{d_c \times d_c} & 0 & \cdots & 0 \\ 0 & \gamma_2 \mathbf{I}_{d_c \times d_c} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \gamma_N \mathbf{I}_{d_c \times d_c} \end{bmatrix}, \quad (17)$$

in which  $\gamma_n \mathbf{I}_{d_c \times d_c}$  is a prior variance  $d_c \times d_c$  matrix of  $\mathbf{x}_n$  and  $\mathbf{I}_{d_c \times d_c}$  is a  $d_c \times d_c$  identity matrix.

The prior distribution  $p(\mathbf{X}|\Gamma^{3D})$  is then defined as

$$p(\mathbf{X}|\Gamma^{3D}) = \prod_{t=1}^T \mathcal{N}(\mathbf{x}(t)|\mathbf{0}, \Gamma^{3D}) = \prod_{t=1}^T \prod_{n=1}^N \mathcal{N}(\mathbf{x}_n(t)|\mathbf{0}, \gamma_n \mathbf{I}_{d_c \times d_c}), \quad (18)$$

As all Type-II algorithms considered here model the source covariance matrix  $\Gamma$  to be diagonal, the proposed extension to 3D sources with free orientation is applicable. The update rule in (14) can be reformulated as follows:

$$\gamma_n^{k+1} \leftarrow \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{\mathbf{x}}_n^k(t))^\top \bar{\mathbf{x}}_n^k(t)}{\text{tr}(\mathbf{L}_n^\top (\Sigma_y^k)^{-1} \mathbf{L}_n)}} \quad \text{for } n = 1, \dots, N. \quad (19)$$

In the simulations studies described in Section IV, we use the fixed-orientation variants of all methods. For real-data analyses in Section V, the free-orientation variants are employed.

**Remark 4.** Note that the number of orientations for the lead field,  $d_c$ , is not limited to 2 or 3. The case  $d_c = 2, 3$  is given when voxel-level lead fields are considered. For lead fields defined for a region or cortical patch level,  $d_c$  can be larger, e.g. determined by the number principal components (PC) describing the voxel lead field within that specific region or path. Interested readers can refer to [88] for more details for such lead field formulations.

#### D. Proof of Theorem 1

*Proof:* We start the proof by recalling (7):

$$\mathcal{L}^{\text{II}}(\Gamma, \Lambda) = -\log p(\mathbf{Y}|\Gamma, \Lambda) = \log |\Sigma_y| + \frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \Sigma_y^{-1} \mathbf{y}(t). \quad (20)$$

The upper bound on the  $\log |\Sigma_y|$  term can be directly inferred from the concavity of the log-determinant function and its first-order Taylor expansion around the value from the previous iteration,  $\Sigma_y^k$ , which provides the following inequality [36, Example 2]:

$$\begin{aligned} \log |\Sigma_y| &\leq \log |\Sigma_y^k| + \text{tr}\left[(\Sigma_y^k)^{-1} (\Sigma_y - \Sigma_y^k)\right] \\ &= \log |\Sigma_y^k| + \text{tr}\left[(\Sigma_y^k)^{-1} \Sigma_y\right] - \text{tr}\left[(\Sigma_y^k)^{-1} \Sigma_y^k\right]. \end{aligned} \quad (21)$$

Note that the first and last terms in (21) do not depend on  $\Gamma$ ; hence, they can be ignored in the optimization procedure. Now, we decompose  $\Sigma_y$  into two terms, each of which only contains either the noise or source covariances:

$$\text{tr}\left[(\Sigma_y^k)^{-1} \Sigma_y\right] = \text{tr}\left[(\Sigma_y^k)^{-1} (\mathbf{L} \mathbf{F} \mathbf{L}^\top + \Lambda)\right] = \text{tr}\left[(\Sigma_y^k)^{-1} \mathbf{L} \mathbf{F} \mathbf{L}^\top\right] + \text{tr}\left[(\Sigma_y^k)^{-1} \Lambda\right]. \quad (22)$$

In next step, we decompose the second term in (7),  $\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \Sigma_y^{-1} \mathbf{y}(t)$ , into two terms, each of which is a function of either only the noise or only the source covariances. To this end, we exploit the following relationship between sensor and source space covariances:

$$\frac{1}{T} \sum_{t=1}^T \mathbf{y}(t)^\top \Sigma_y^{-1} \mathbf{y}(t) = \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t)^\top \Gamma^{-1} \bar{\mathbf{x}}^k(t) + \frac{1}{T} \sum_{t=1}^T (\mathbf{y}(t) - \mathbf{L} \bar{\mathbf{x}}^k(t))^\top \Lambda^{-1} (\mathbf{y}(t) - \mathbf{L} \bar{\mathbf{x}}^k(t)). \quad (23)$$

By combining (22) and (23), rearranging the terms, and ignoring all terms that do not depend on  $\Gamma$ , we have:

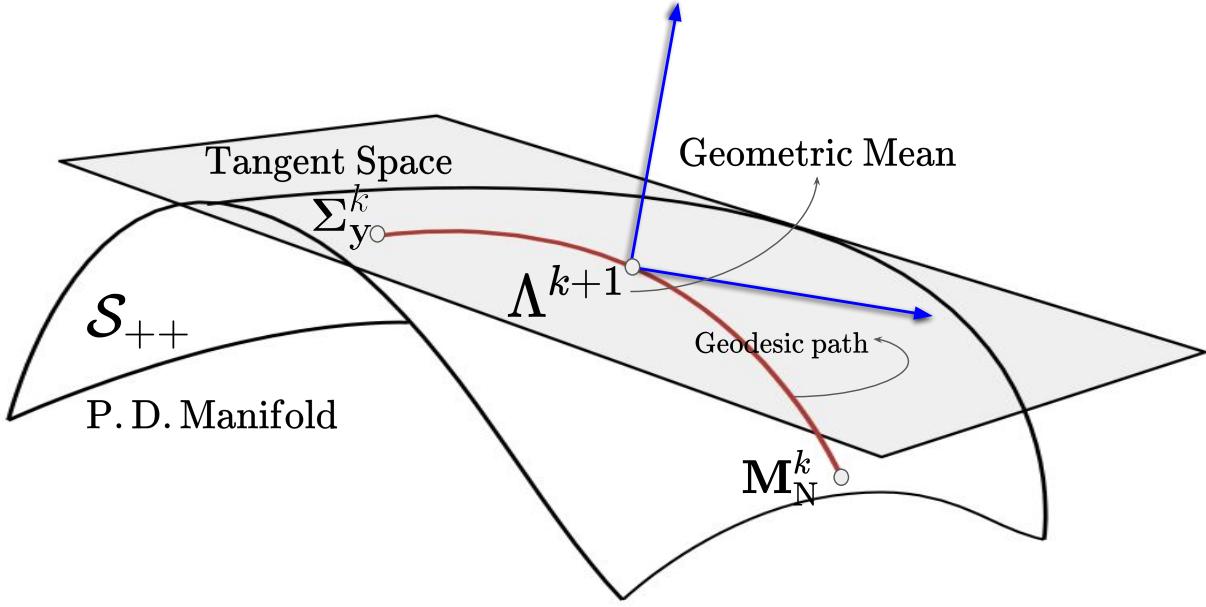
$$\begin{aligned} \mathcal{L}^{\text{II}}(\Gamma) &\leq \text{tr}\left[(\Sigma_y^k)^{-1} \mathbf{L} \mathbf{F} \mathbf{L}^\top\right] + \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t)^\top \Gamma^{-1} \bar{\mathbf{x}}^k(t) + \text{const} \\ &= \text{tr}\left(\mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L} \Gamma\right) + \text{tr}(\mathbf{M}_S^k \Gamma^{-1}) + \text{const} = \mathcal{L}_{\text{source}}^{\text{conv}}(\Gamma, \Lambda^k) + \text{const}, \end{aligned} \quad (24)$$

where  $\mathbf{M}_S^k := \frac{1}{T} \sum_{t=1}^T \bar{\mathbf{x}}^k(t) \bar{\mathbf{x}}^k(t)^\top$ .

Note that constant values in (24) do not depend on  $\Gamma$ ; hence, they can be ignored in the optimization procedure. This proves the equivalence of (7) and (8) when the optimization is performed with respect to  $\Gamma$ .

The equivalence of (7) and (10) can be shown analogously, with the difference that we only focus on noise-related terms in (22) and (23):

$$\begin{aligned} \mathcal{L}^{\text{II}}(\Lambda) &\leq \text{tr}\left[(\Sigma_y^k)^{-1} \Lambda\right] + \frac{1}{T} \sum_{t=1}^T (\mathbf{y}(t) - \mathbf{L} \bar{\mathbf{x}}^k(t))^\top \Lambda^{-1} (\mathbf{y}(t) - \mathbf{L} \bar{\mathbf{x}}^k(t)) + \text{const} \\ &= \text{tr}\left[(\Sigma_y^k)^{-1} \Lambda\right] + \text{tr}(\mathbf{M}_N^k \Lambda^{-1}) + \text{const} = \mathcal{L}_{\text{noise}}^{\text{conv}}(\Gamma^k, \Lambda) + \text{const}, \end{aligned} \quad (25)$$



**Fig. 7:** Geometric representation of the geodesic path between the pair of matrices  $\{\Sigma_y^k, M_N^k\}$  on the P.D. manifold and the geometric mean between them, which is used to update  $\Lambda^{k+1}$ .

where  $M_N^k := \frac{1}{T} \sum_{t=1}^T (\mathbf{y}(t) - \mathbf{L}\bar{x}^k(t))(\mathbf{y}(t) - \mathbf{L}\bar{x}^k(t))^\top$ .

Constant values in (25) do not depend on  $\Lambda$ ; hence, they can again be ignored in the optimization procedure. Summarizing, we have shown that optimizing (7) is equivalent to optimizing  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\Gamma^k, \Lambda)$  and  $\mathcal{L}_{\text{source}}^{\text{conv}}(\Gamma, \Lambda^k)$ , which concludes the proof.

### E. 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  $\{\Sigma_y^k, M_N^k\}$ .

**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 [89, Chapter. 6]:*

$$\xi_l(\mathbf{P}_0, \mathbf{P}_1) = (\mathbf{P}_0)^{\frac{1}{2}} \left( (\mathbf{P}_0)^{-1/2} \mathbf{P}_1 (\mathbf{P}_0)^{-1/2} \right)^l (\mathbf{P}_0)^{\frac{1}{2}} \quad l \in [0, 1]. \quad (26)$$

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 :

$$\begin{aligned} \xi_l(\Lambda_0, \Lambda_1)^{-1} &= \left( (\Lambda_0)^{1/2} \left( (\Lambda_0)^{-1/2} \Lambda_1 (\Lambda_0)^{-1/2} \right)^l (\Lambda_0)^{1/2} \right)^{-1} \\ &= \left( (\Lambda_0)^{-1/2} \left( (\Lambda_0)^{1/2} (\Lambda_1)^{-1} (\Lambda_0)^{1/2} \right)^l (\Lambda_0)^{-1/2} \right) = \xi_l(\Lambda_0^{-1}, \Lambda_1^{-1}). \end{aligned} \quad (27)$$

**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  $f$  with domain  $\mathcal{A} \subset \mathcal{M}$  with  $f : \mathcal{A} \rightarrow \mathbb{R}$  is called geodesic convex (*g-convex*) if the following relation holds:

$$f(\zeta_l(\mathbf{p}_0, \mathbf{p}_1)) \leq lf(\mathbf{p}_0) + (1-l)f(\mathbf{p}_1), \quad (28)$$

where  $l \in [0, 1]$  and  $\zeta_l(\mathbf{p}_0, \mathbf{p}_1)$  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  $\zeta_l(\mathbf{p}_0, \mathbf{p}_1)$  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 5.** 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:

$$\zeta_l(\mathbf{p}_0, \mathbf{p}_1) = l\mathbf{p}_0 + (1-l)\mathbf{p}_1. \quad (29)$$

For the sake of brevity, we omit a detailed theoretical introduction of *g-convexity*, and only borrow a result from [90]. Interested readers are referred to [91, Chapter 1] for a gentle introduction to this topic, and [92, Chapter. 2]; [85], [93]–[101] for more in-depth technical details. Now we are ready to state the proof, which parallels the one provided in [90, Theorem. 3].

*Proof:* We proceed in two steps. First, we consider P.D. manifolds and express (28) in terms of geodesic paths and functions that lie on this particular space. We then show that  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda})$  is strictly *g-convex* on this specific domain. In the second step, we then derive the update rule proposed in (12).

**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  $\xi_l(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)$  as presented in Definition 2, and define  $f(\cdot)$  to be  $f(\boldsymbol{\Lambda}) = \text{tr}\left[\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\boldsymbol{\Lambda}\right] + \text{tr}(\mathbf{M}_N^k\boldsymbol{\Lambda}^{-1})$ , representing the cost function  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda})$ .

We now show that  $f(\boldsymbol{\Lambda})$  is strictly *g-convex* on this specific domain. For continuous functions as considered in this paper, fulfilling (28) for  $f(\boldsymbol{\Lambda})$  and  $\xi_l(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)$  with  $l = 1/2$  is sufficient for strict *g-convexity* according to *mid-point convexity* [102]:

$$\begin{aligned} & \text{tr}\left(\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\xi_{1/2}(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)\right) + \text{tr}\left(\mathbf{M}_N^k\xi_{1/2}(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)^{-1}\right) \\ & < \frac{1}{2}\text{tr}\left(\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\boldsymbol{\Lambda}_0\right) + \frac{1}{2}\text{tr}\left(\mathbf{M}_N^k\boldsymbol{\Lambda}_0^{-1}\right) \\ & \quad + \frac{1}{2}\text{tr}\left(\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\boldsymbol{\Lambda}_1\right) + \frac{1}{2}\text{tr}\left(\mathbf{M}_N^k\boldsymbol{\Lambda}_1^{-1}\right). \end{aligned} \quad (30)$$

Given  $\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1} \in \mathcal{S}_{++}$ , i.e.,  $\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1} > 0$  and the operator inequality [89, Chapter. 4]

$$\xi_{1/2}(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1) \prec \frac{1}{2}\boldsymbol{\Lambda}_0 + \frac{1}{2}\boldsymbol{\Lambda}_1, \quad (31)$$

we have:

$$\text{tr}\left(\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\xi_{1/2}(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)\right) < \frac{1}{2}\text{tr}\left(\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\boldsymbol{\Lambda}_0\right) + \frac{1}{2}\text{tr}\left(\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\boldsymbol{\Lambda}_1\right), \quad (32)$$

which is derived by multiplying both sides of (31) with  $\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}$  followed by taking the trace on both sides.

Similarly, we can write the operator inequality for  $\{\boldsymbol{\Lambda}_0^{-1}, \boldsymbol{\Lambda}_1^{-1}\}$  using (27) as:

$$\xi_{1/2}(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)^{-1} = \xi_{1/2}(\boldsymbol{\Lambda}_0^{-1}, \boldsymbol{\Lambda}_1^{-1}) \prec \frac{1}{2}\boldsymbol{\Lambda}_0^{-1} + \frac{1}{2}\boldsymbol{\Lambda}_1^{-1}. \quad (33)$$

Multiplying both sides of (33) by  $\mathbf{M}_N^k \in \mathcal{S}_{++}$  and applying the trace operator on both sides leads to:

$$\text{tr}\left(\mathbf{M}_N^k\xi_{1/2}(\boldsymbol{\Lambda}_0, \boldsymbol{\Lambda}_1)^{-1}\right) < \frac{1}{2}\text{tr}\left(\mathbf{M}_N^k\boldsymbol{\Lambda}_0^{-1}\right) + \frac{1}{2}\text{tr}\left(\mathbf{M}_N^k\boldsymbol{\Lambda}_1^{-1}\right). \quad (34)$$

Summing up (32) and (34) proves inequality (30) and concludes the first part of the proof.

**2) Part II: Derivation of the Update Rule in (12):** We now present the second part of the proof by deriving the update rule in (12). Since the cost function  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda})$  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 (10) and setting the result to zero as follows:

$$\nabla\mathcal{L}_{\text{noise}}^{\text{conv}}(\boldsymbol{\Gamma}^k, \boldsymbol{\Lambda}) = \left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1} - \boldsymbol{\Lambda}^{-1}\mathbf{M}_N^k\boldsymbol{\Lambda}^{-1} = 0, \quad (35)$$

which results in

$$\boldsymbol{\Lambda}\left(\boldsymbol{\Sigma}_{\mathbf{y}}^k\right)^{-1}\boldsymbol{\Lambda} = \mathbf{M}_N^k. \quad (36)$$

This solution is known as the *Riccati equation*, and is the geometric mean between  $\Sigma_y^k$  and  $M_N^k$  [103], [104]:

$$\Lambda^{k+1} \leftarrow (\Sigma_y^k)^{\frac{1}{2}} \left( (\Sigma_y^k)^{-1/2} M_N^k (\Sigma_y^k)^{-1/2} \right)^{\frac{1}{2}} (\Sigma_y^k)^{\frac{1}{2}}.$$

A geometric representation of the geodesic path between the pair of matrices  $\{\Sigma_y^k, M_N^k\}$  on the P.D. manifold and the geometric mean between them, representing the update for  $\Lambda^{k+1}$ , is provided in Fig. 7.

**Remark 6.** Note that the obtained update rule is a closed-form solution for the surrogate cost function, (10), which stands in contrast to conventional majorization minimization algorithms (see section G in the appendix), which require iterative procedures in each step of the optimization.

Deriving the update rule in (12) concludes the second part of the proof of Theorem 2.

#### F. Proof of Theorem 3

We start the derivation of update rule (14) by constraining  $\Gamma$  to the set of diagonal matrices with non-negative entries  $\mathcal{S}$ , i.e.,

$$\mathcal{S} = \{\Gamma \mid \Gamma = \text{diag}(\gamma) = \text{diag}([\gamma_1, \dots, \gamma_N]^\top), \gamma_n \geq 0, \text{ for } n = 1, \dots, N\}.$$

We continue by reformulating the constrained optimization with respect to the source covariance matrix,

$$\Gamma^{k+1} = \underset{\Gamma \in \mathcal{S}, \Lambda = \Lambda^k}{\arg \min} \text{tr} \left( \mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L} \Gamma \right) + \text{tr}(M_S^k \Gamma^{-1}), \quad (37)$$

as follows:

$$\gamma^{k+1} = \underset{\gamma \geq 0, \Lambda = \Lambda^k}{\arg \min} \underbrace{\text{diag} \left[ \mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L} \right] \gamma + \text{diag} [M_S^k] \gamma^{-1}}_{\mathcal{L}_{\text{source}}^{\text{diag}}(\gamma | \gamma^k)}, \quad (38)$$

where  $\gamma^{-1} = [\gamma_1^{-1}, \dots, \gamma_N^{-1}]^\top$  is defined as the element-wise inversion of  $\gamma$ . Note that the set of diagonal matrices with all non-negative entries are positive semidefinite (PSD) by construction [105, Appendix A]. Thus, by constraining the space of solutions of optimization problem (37) to the set  $\mathcal{S}$ , the PSD requirement for  $\Gamma$  reduces to the requirement that the diagonal elements of  $\Gamma$ , i.e.,  $\gamma_n$ , for  $n = 1, \dots, N$ , must be non-negative. The optimization with respect to the scalar source variances is then carried out by taking the derivative of (38) with respect to  $\gamma_n$ , for  $n = 1, \dots, N$ , and setting it to zero,

$$\begin{aligned} \frac{\partial}{\partial \gamma_n} & \left( \left[ \mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L} \right] \gamma_n + [M_S^k] \gamma_n^{-1} \right) \\ &= \left[ \mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L} \right]_{n,n} - \frac{1}{(\gamma_n)^2} [M_S^k]_{n,n} \\ &= 0 \quad \text{for } n = 1, \dots, N, \end{aligned}$$

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

$$\Gamma^{k+1} = \text{diag}(\gamma^{k+1}), \text{ where, } \gamma_n^{k+1} \leftarrow \sqrt{\frac{[M_S^k]_{n,n}}{\left[ \mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L} \right]_{n,n}}} = \sqrt{\frac{\frac{1}{T} \sum_{t=1}^T (\bar{x}_n^k(t))^2}{\mathbf{L}_n^\top (\Sigma_y^k)^{-1} \mathbf{L}_n}} \quad \text{for } n = 1, \dots, N,$$

which is identical to the update rule of Champagne [30].

#### G. Proof of Theorem 4

We prove Theorem 4 by showing that the alternating update rules for  $\Lambda$  and  $\Gamma$ , (12) and (14), are guaranteed to converge to a local minimum of the Bayesian Type-II likelihood (7). More generally, we prove that FUN learning 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 [60], [61] [59], [106].

**1) Required Conditions for Majorization-Minimization Algorithms:** MM encompasses a family of iterative algorithms for optimizing general non-linear cost functions. 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. The MM class covers a broad range of common optimization algorithms such as *convex-concave procedures (CCCP)* and *proximal methods* [36, Section IV], [108]–[110]. Such algorithms have been applied in various domains such as brain source imaging [111], [112] [6], [7], wireless communication systems with massive MIMO technology [72], [113], [114], and non-negative matrix factorization [115]. Interested readers are referred to [36] for an extensive list of applications on MM.

The problem of minimizing a continuous function  $f(\mathbf{u})$  within a closed convex set  $\mathcal{U} \subset \mathbb{R}^n$ :

$$\min_{\mathbf{u}} f(\mathbf{u}) \quad \text{subject to } \mathbf{u} \in \mathcal{U}, \quad (39)$$

within the MM framework 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

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

If a surrogate function fulfills conditions [A1]–[A3], then the value of the cost function  $f$  decreases in each iteration:  $f(\mathbf{u}^{k+1}) \leq f(\mathbf{u}^k)$ . For the smooth functions considered in this paper, we further 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 theorem:

**Theorem 5.** *Assume that an MM algorithm fulfills conditions [A1]–[A4]. Then, every limit point of the sequence of minimizers generated in [A3], is a stationary point of the original optimization problem in (39).*

*Proof:* A detailed proof is provided in [61, Theorem 1].

**2) Detail Derivation of the Proof of Theorem 4:** We now show that FUN learning is an instance of majorization-minimization as defined above, which fulfills Theorem 5.

*Proof:* We need to prove that conditions [A1]–[A4] are fulfilled for FUN learning. To this end, we recall the upper bound on  $\log |\Sigma_y|$  in (21), which fulfills condition [A2] since it majorizes  $\log |\Sigma_y|$  by virtue of the concavity of the log-determinant function and its first-order Taylor expansion around  $\Sigma_y^k$ . Besides, it automatically satisfies conditions [A1] and [A4] by construction, because the majorizing function in (21) is obtained through a Taylor expansion around  $\Sigma_y^k$ . Concretely, [A1] is satisfied because the equality in (21) holds for  $\Sigma_y = \Sigma_y^k$ . Similarly, [A4] is satisfied because the gradient of  $\log |\Sigma_y|$  at point  $\Sigma_y^k$ ,  $(\Sigma_y^k)^{-1}$  defines the linear Taylor approximation  $\log |\Sigma_y^k| + \text{tr}[(\Sigma_y^k)^{-1} (\Sigma_y - \Sigma_y^k)]$ . Thus, both gradients coincide in  $\Sigma_y^k$  by construction. We can further prove that [A3] can be satisfied by showing that  $\hat{\mathcal{L}}_{\text{noise}}^{\text{conv}}(\Gamma^k, \Lambda)$  reaches its global minimum in each MM iteration. This is guaranteed if  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\Gamma^k, \Lambda)$  can be shown to be convex or g-convex with respect to  $\Lambda$ . To this end, we first require the subsequent proposition:

**Proposition 1.** *Any local minimum of a g-convex function over a g-convex set is a global minimum.*

*Proof:* A detailed proof is presented in [93, Theorem 2.1].

Given the proof presented in appendix E.1, we can conclude that  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\Gamma^k, \Lambda)$  is g-convex; hence, any local minimum of  $\mathcal{L}_{\text{noise}}^{\text{conv}}(\Gamma^k, \Lambda)$  is a global minimum according to Proposition 1. This proves that condition [A3] is fulfilled and completes the proof that the optimization of (7) with respect to  $\Lambda$  using the convex surrogate cost function (10) leads to an MM algorithm.

We omit the proof of conditions [A1], [A2] and [A4] for the optimization with respect to  $\Gamma$  based on the convex surrogate function in (8),  $\mathcal{L}_{\text{source}}^{\text{conv}}(\Gamma, \Lambda^k)$ , as it can be presented analogously. We here only show that [A3] is satisfied if  $\mathcal{L}_{\text{source}}^{\text{diag}}(\gamma|\gamma^k)$  in (38) is a convex function with respect to  $\gamma$ . Note that the g-convexity of  $\mathcal{L}_{\text{source}}^{\text{conv}}(\Gamma, \Lambda^k)$  can also be proven using arguments analogous to those presented in appendix E.1. However, we instead prove a stronger condition, i.e., convexity, for simplifying the proof. To this end, we rewrite (38) as follows:

$$\mathcal{L}_{\text{source}}^{\text{diag}}(\gamma|\gamma^k) = \text{diag}[\mathbf{V}^k] \gamma + \text{diag}[\mathbf{M}_S^k] \gamma^{-1},$$

where  $\mathbf{V}^k := \mathbf{L}^\top (\Sigma_y^k)^{-1} \mathbf{L}$  is defined as a parameter that does not depend on  $\gamma$ . The convexity of  $\mathcal{L}_{\text{source}}^{\text{diag}}(\gamma|\gamma^k)$  can be directly inferred from the convexity of  $\text{diag}[\mathbf{V}^k] \gamma$  and  $\text{diag}[\mathbf{M}_S^k] \gamma^{-1}$  with respect to  $\gamma$  [116, Chapter. 3]. The convexity of  $\mathcal{L}_{\text{source}}^{\text{diag}}(\gamma|\gamma^k)$ , which ensures that condition [A3] can be satisfied using standard optimization, along with the fulfillment of

conditions [A1], [A2] and [A4], ensure that Theorem 5 holds for  $\mathcal{L}_{\text{source}}^{\text{conv}}(\boldsymbol{\Gamma}, \boldsymbol{\Lambda}^k)$ . This completes the proof that the optimization of (7) with respect to  $\boldsymbol{\Gamma}$  using the convex surrogate cost function (8) leads to an MM algorithm with convergence guarantees.

#### H. Derivation of Champagne with Heteroscedastic Noise Learning as a Special Case of FUN Learning

Similar to Appendix F, we start by constraining  $\boldsymbol{\Lambda}$  to the set of diagonal matrices with non-negative entries  $\mathcal{S}$ , i.e.,

$$\mathcal{S} = \{\boldsymbol{\Lambda} \mid \boldsymbol{\Lambda} = \text{diag}(\boldsymbol{\lambda}) = \text{diag}([\lambda_1, \dots, \lambda_M]^T), \lambda_m \geq 0, \text{ for } m = 1, \dots, M\}.$$

We continue by reformulating the constrained optimization with respect to the noise covariance matrix,

$$\boldsymbol{\Lambda}^{k+1} = \underset{\boldsymbol{\Lambda} \in \mathcal{S}, \boldsymbol{\Gamma} = \boldsymbol{\Gamma}^k}{\arg \min} \text{tr} \left( (\boldsymbol{\Sigma}_y^k)^{-1} \boldsymbol{\Lambda} \right) + \text{tr}(\mathbf{M}_N^k \boldsymbol{\Lambda}^{-1}), \quad (40)$$

as follows:

$$\boldsymbol{\lambda}^{k+1} = \underset{\boldsymbol{\lambda} \geq 0, \boldsymbol{\Gamma} = \boldsymbol{\Gamma}^k}{\arg \min} \underbrace{\text{diag} \left[ (\boldsymbol{\Sigma}_y^k)^{-1} \right] \boldsymbol{\lambda} + \text{diag} \left[ \mathbf{M}_N^k \right] \boldsymbol{\lambda}^{-1}}_{\mathcal{L}_{\text{noise}}^{\text{diag}}(\boldsymbol{\lambda} | \boldsymbol{\lambda}^k)}, \quad (41)$$

where  $\boldsymbol{\lambda}^{-1} = [\lambda_1^{-1}, \dots, \lambda_M^{-1}]^T$  is defined as the element-wise inversion of  $\boldsymbol{\lambda}$ . Taking the derivative of (41) with respect to  $\lambda_m$ , for  $m = 1, \dots, M$ , and setting it to zero,

$$\begin{aligned} \frac{\partial}{\partial \lambda_m} & \left( \left[ (\boldsymbol{\Sigma}_y^k)^{-1} \right] \lambda_m + \left[ \mathbf{M}_N^k \right] \lambda_m^{-1} \right) \\ &= \left[ (\boldsymbol{\Sigma}_y^k)^{-1} \right]_{m,m} - \frac{1}{(\lambda_m)^2} \left[ \mathbf{M}_N^k \right]_{m,m} \\ &= 0 \quad \text{for } m = 1, \dots, M, \end{aligned}$$

yields the following update rule:

$$\lambda_m^{k+1} \leftarrow \sqrt{\frac{\left[ \mathbf{M}_N^k \right]_{m,m}}{\left[ (\boldsymbol{\Sigma}_y^k)^{-1} \right]_{m,m}}} = \sqrt{\frac{\left[ \frac{1}{T} \sum_{t=1}^T (\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t))(\mathbf{y}(t) - \mathbf{L}\bar{\mathbf{x}}^k(t))^T \right]_{m,m}}{\left[ (\boldsymbol{\Sigma}_y^k)^{-1} \right]_{m,m}}} \quad \text{for } m = 1, \dots, M, \quad (42)$$

which is identical to the update rule of the Champagne with heteroscedastic noise learning as presented in [7].

#### ACKNOWLEDGMENT

This result is part of a project that has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (Grant agreement No. 758985).

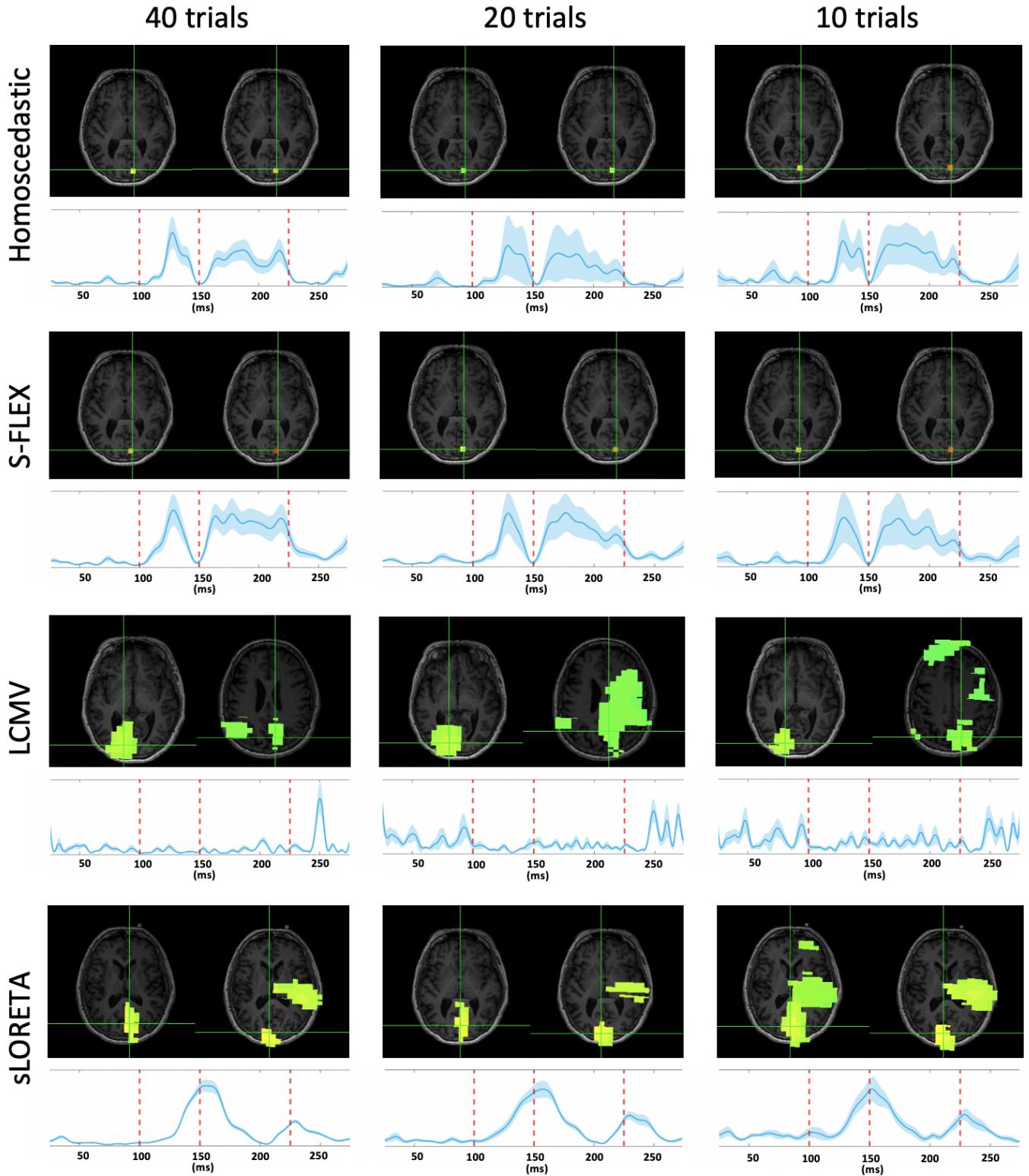
AH acknowledges scholarship support from the Machine Learning/Intelligent Data Analysis research group at Technische Universität Berlin. He further wishes to thank the Charité – Universitätsmedizin Berlin, the Berlin Mathematical School (BMS), and the Berlin Mathematics Research Center MATH+ for partial support. CC was supported by Hubei Provincial Natural Science Foundation of China under Grant 2021CFB384 and the National Natural Science Foundation of China under Grant 62007013. KRM was partly funded by the German Ministry for Education and Research (under refs 01IS14013A-E, 01GQ1115, 01GQ0850, 01IS18056A, 01IS18025A and 01IS18037A), the German Research Foundation (DFG) as Math+: Berlin Mathematics Research Center (EXC 2046/1, project-ID: 390685689). Furthermore, KRM was partly supported by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grants funded by the Korea Government (No. 2019-0-00079, Artificial Intelligence Graduate School Program, Korea University). SSN was funded in part by National Institutes of Health grants (R01DC004855, R01EB022717, R01DC176960, R01DC010145, R01NS100440, R01AG062196, and R01DC013979), University of California MRPI MRP-17-454755, the US Department of Defense grant (W81XWH-13-1-0494).

#### REFERENCES

- [66] A. Wu, O. Koyejo, and J. W. Pillow, "Dependent relevance determination for smooth and structured sparse regression." *J. Mach. Learn. Res.*, vol. 20, no. 89, pp. 1–43, 2019.
- [67] O. Dikmen and C. Févotte, "Maximum marginal likelihood estimation for nonnegative dictionary learning in the Gamma-Poisson model," *IEEE Transactions on Signal Processing*, vol. 60, no. 10, pp. 5163–5175, 2012.
- [68] F. R. Bach and M. I. Jordan, "A probabilistic interpretation of canonical correlation analysis," *Technical Report 688, Department of Statistics, University of California, Berkeley*, 2005.
- [69] P. Ablin, J.-F. Cardoso, and A. Gramfort, "Spectral independent component analysis with noise modeling for M/EEG source separation," *arXiv preprint arXiv:2008.09693*, 2020.
- [70] R. Prasad, C. R. Murthy, and B. D. Rao, "Joint channel estimation and data detection in MIMO-OFDM systems: A sparse Bayesian learning approach," *IEEE Transactions on Signal Processing*, vol. 63, no. 20, pp. 5369–5382, 2015.

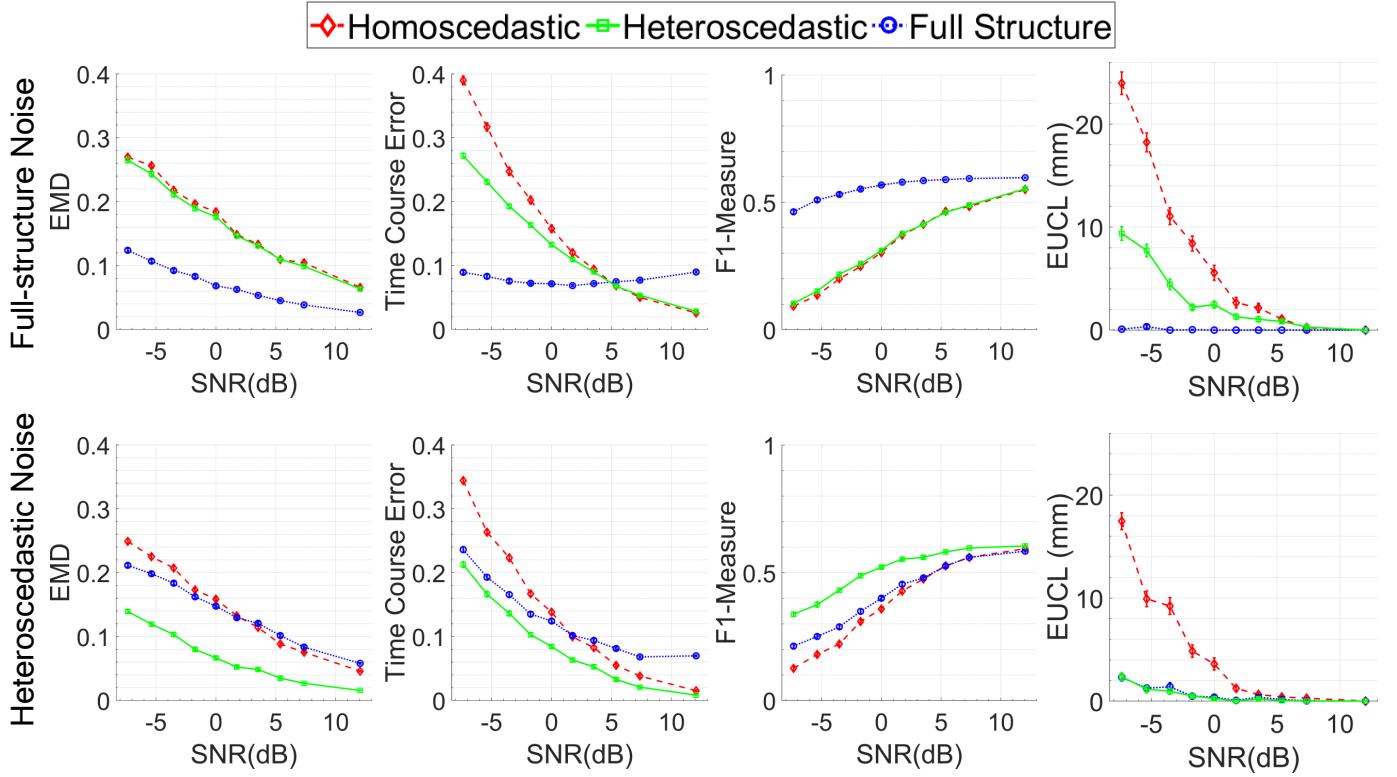
- [71] P. Gerstoft, C. F. Mecklenbräuker, A. Xenaki, and S. Nannuru, "Multisnapshot sparse Bayesian learning for DOA," *IEEE Signal Processing Letters*, vol. 23, no. 10, pp. 1469–1473, 2016.
- [72] S. Haghighatshoar and G. Caire, "Massive MIMO channel subspace estimation from low-dimensional projections," *IEEE Transactions on Signal Processing*, vol. 65, no. 2, pp. 303–318, 2017.
- [73] Y. Feng, D. P. Palomar *et al.*, "A signal processing perspective on financial engineering," *Foundations and Trends® in Signal Processing*, vol. 9, no. 1–2, pp. 1–231, 2016.
- [74] K. Werner, M. Jansson, and P. Stoica, "On estimation of covariance matrices with Kronecker product structure," *IEEE Transactions on Signal Processing*, vol. 56, no. 2, pp. 478–491, 2008.
- [75] T. Tsiligkaridis and A. O. Hero, "Covariance estimation in high dimensions via Kronecker product expansions," *IEEE Transactions on Signal Processing*, vol. 61, no. 21, pp. 5347–5360, 2013.
- [76] A. M. Zoubir, V. Koivunen, E. Ollila, and M. Muma, *Robust statistics for signal processing*. Cambridge University Press, 2018.
- [77] A. Benfenati, E. Chouzenoux, and J.-C. Pesquet, "Proximal approaches for matrix optimization problems: Application to robust precision matrix estimation," *Signal Processing*, vol. 169, p. 107417, 2020.
- [78] E. Ollila, D. P. Palomar, and F. Pascal, "Shrinking the eigenvalues of M-estimators of covariance matrix," *arXiv preprint arXiv:2006.10005*, 2020.
- [79] F. Bouchard, A. Breloy, G. Ginolhac, A. Renaux, and F. Pascal, "A riemannian framework for low-rank structured elliptical models," *arXiv preprint arXiv:2001.01141*, 2020.
- [80] B. Meriaux, C. Ren, A. Breloy, M. N. El Korso, and P. Forster, "Matched and mismatched estimation of kronecker product of linearly structured scatter matrices under elliptical distributions," *IEEE Transactions on Signal Processing*, 2020.
- [81] S. Kumar, J. Ying, J. V. de Miranda Cardoso, and D. P. Palomar, "A unified framework for structured graph learning via spectral constraints," *Journal of Machine Learning Research*, vol. 21, no. 22, pp. 1–60, 2020.
- [82] A. Hashemi, M. Rostami, and N.-M. Cheung, "Efficient environmental temperature monitoring using compressed sensing," in *2016 Data Compression Conference (DCC)*. IEEE, 2016, pp. 602–602.
- [83] A. Flinth and A. Hashemi, "Thermal source localization through infinite-dimensional compressed sensing," *arXiv preprint arXiv:1710.02016*, 2017.
- [84] ———, "Approximate recovery of initial point-like and instantaneous sources from coarsely sampled thermal fields via infinite-dimensional compressed sensing," in *2018 26th European Signal Processing Conference (EUSIPCO)*. IEEE, 2018, pp. 1720–1724.
- [85] H. Wei, A. Jafarian, P. Zeidman, V. Litvak, A. Razi, D. Hu, and K. J. Friston, "Bayesian fusion and multimodal DCM for EEG and fMRI," *NeuroImage*, vol. 211, p. 116595, 2020.
- [86] R. Giordano, T. Broderick, and M. I. Jordan, "Covariances, robustness and variational Bayes," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 1981–2029, 2018.
- [87] B. Dai, Y. Wang, J. Aston, G. Hua, and D. Wipf, "Connections with robust PCA and the role of emergent sparsity in variational autoencoder models," *The Journal of Machine Learning Research*, vol. 19, no. 1, pp. 1573–1614, 2018.
- [88] T. Limpiti, B. D. Van Veen, and R. T. Wakai, "Cortical patch basis model for spatially extended neural activity," *IEEE Transactions on Biomedical Engineering*, vol. 53, no. 9, pp. 1740–1754, 2006.
- [89] R. Bhatia, *Positive definite matrices*. Princeton University Press, 2009, vol. 24.
- [90] P. Zadeh, R. Hosseini, and S. Sra, "Geometric mean metric learning," in *International Conference on Machine Learning*, 2016, pp. 2464–2471.
- [91] A. Wiesel, T. Zhang *et al.*, "Structured robust covariance estimation," *Foundations and Trends® in Signal Processing*, vol. 8, no. 3, pp. 127–216, 2015.
- [92] A. Papadopoulos, *Metric spaces, convexity and nonpositive curvature*. European Mathematical Society, 2005, vol. 6.
- [93] T. Rapcsak, "Geodesic convexity in nonlinear optimization," *Journal of Optimization Theory and Applications*, vol. 69, no. 1, pp. 169–183, 1991.
- [94] A. Ben-Tal, "On generalized means and generalized convex functions," *Journal of Optimization Theory and Applications*, vol. 21, no. 1, pp. 1–13, 1977.
- [95] L. Liberti, "On a class of nonconvex problems where all local minima are global," *Publications de l'Institut Mathématique*, vol. 76, no. 90, pp. 101–109, 2004.
- [96] D. E. Pallaschke and S. Rolewicz, *Foundations of mathematical optimization: convex analysis without linearity*. Springer Science & Business Media, 2013, vol. 388.
- [97] M. Moakher, "A differential geometric approach to the geometric mean of symmetric positive-definite matrices," *SIAM Journal on Matrix Analysis and Applications*, vol. 26, no. 3, pp. 735–747, 2005.
- [98] S. Sra and R. Hosseini, "Conic geometric optimization on the manifold of positive definite matrices," *SIAM Journal on Optimization*, vol. 25, no. 1, pp. 713–739, 2015.
- [99] N. K. Vishnoi, "Geodesic convex optimization: Differentiation on manifolds, geodesics, and convexity," *arXiv preprint arXiv:1806.06373*, 2018.
- [100] M. Berger, *A panoramic view of Riemannian geometry*. Springer Science & Business Media, 2012.
- [101] J. Jost, *Riemannian geometry and geometric analysis*. Springer, 2008, vol. 42005.
- [102] C. Niculescu and L.-E. Persson, *Convex functions and their applications*. Springer, 2006.
- [103] J. V. Davis, B. Kulis, P. Jain, S. Sra, and I. S. Dhillon, "Information-theoretic metric learning," in *Proceedings of the 24th International Conference on Machine Learning*, 2007, pp. 209–216.
- [104] S. Bonnabel and R. Sepulchre, "Riemannian metric and geometric mean for positive semidefinite matrices of fixed rank," *SIAM Journal on Matrix Analysis and Applications*, vol. 31, no. 3, pp. 1055–1070, 2009.
- [105] E. De Klerk, *Aspects of semidefinite programming: interior point algorithms and selected applications*. Springer Science & Business Media, 2006, vol. 65.
- [106] D. R. Hunter and K. Lange, "A tutorial on MM algorithms," *The American Statistician*, vol. 58, no. 1, pp. 30–37, 2004.
- [107] T. T. Wu, K. Lange *et al.*, "The MM alternative to EM," *Statistical Science*, vol. 25, no. 4, pp. 492–505, 2010.
- [108] E. Mjolsness and C. Garrett, "Algebraic transformations of objective functions," *Neural Networks*, vol. 3, no. 6, pp. 651–669, 1990.
- [109] A. L. Yuille and A. Rangarajan, "The concave-convex procedure," *Neural computation*, vol. 15, no. 4, pp. 915–936, 2003.
- [110] T. Lipp and S. Boyd, "Variations and extension of the convex–concave procedure," *Optimization and Engineering*, vol. 17, no. 2, pp. 263–287, 2016.
- [111] A. Hashemi and S. Haufe, "Improving EEG source localization through spatio-temporal sparse Bayesian learning," in *2018 26th European Signal Processing Conference (EUSIPCO)*. IEEE, 2018, pp. 1935–1939.
- [112] Y. Bekhti, F. Lucka, J. Salmon, and A. Gramfort, "A hierarchical Bayesian perspective on majorization-minimization for non-convex sparse regression: application to M/EEG source imaging," *Inverse Problems*, vol. 34, no. 8, p. 085010, 2018.
- [113] M. Masood, A. Ghrayeb, P. Babu, I. Khalil, and M. Hasna, "A minorization-maximization algorithm for an-based MIMOE secrecy rate maximization," in *2016 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2016, pp. 975–980.
- [114] M. B. Khalilsarai, T. Yang, S. Haghighatshoar, and G. Caire, "Structured channel covariance estimation from limited samples in Massive MIMO," in *ICC 2020-2020 IEEE International Conference on Communications (ICC)*. IEEE, 2020, pp. 1–7.
- [115] D. Fagot, H. Wendt, C. Févotte, and P. Smaragdis, "Majorization-minimization algorithms for convolutive NMF with the beta-divergence," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 8202–8206.
- [116] S. P. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.

## SUPPLEMENTARY FIGURE

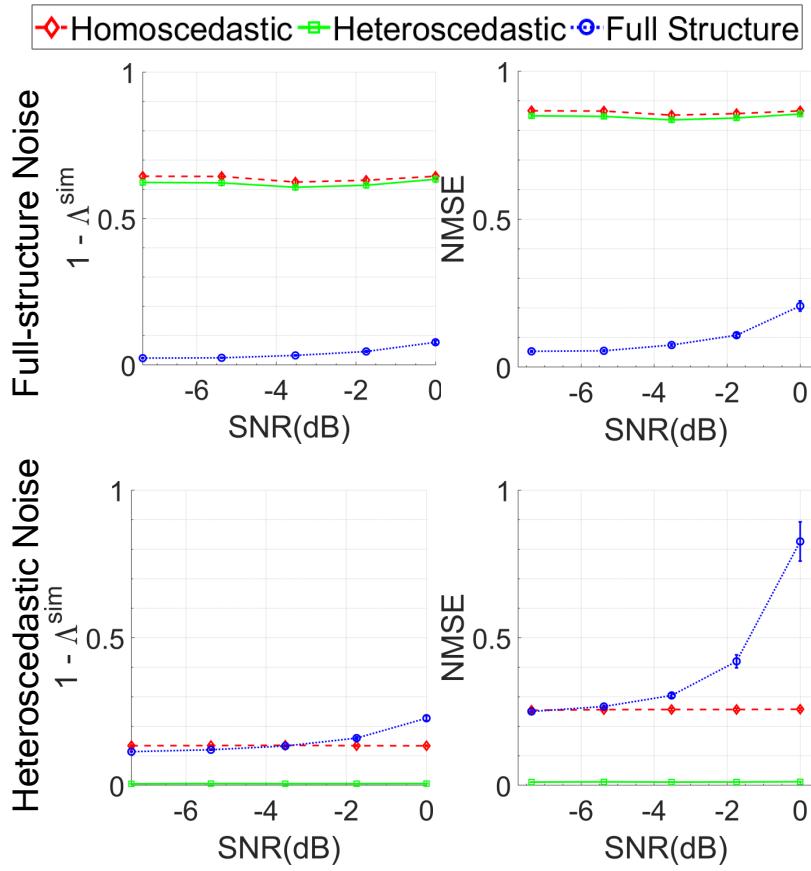


**Fig. 8:** Localization and time series results of visual evoked field (VEF) activity for a single subject using another four benchmark algorithms. As is shown, FUN outperforms LCMV beamformer and sLORETA in terms of localization. Moreover, the activation time courses derived from homoscedastic noise learning Champagne and S-FLEX do not exhibit sharp responses as observed for FUN. The noise level used for S-FLEX reconstructions was set to values learnt from classical Champagne algorithm with noise learning.

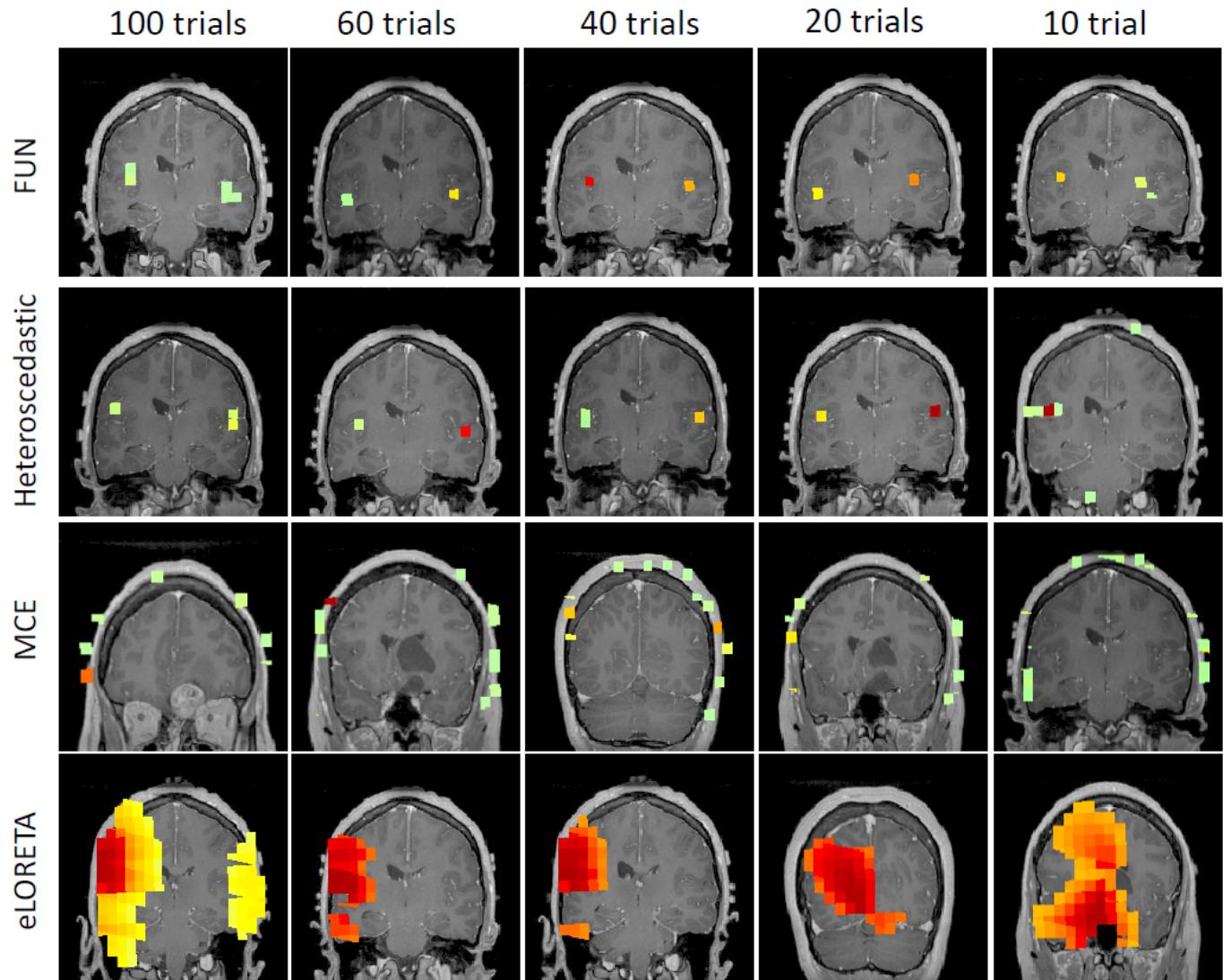
## MAIN FIGURES OF THE PAPER WITH HIGHER RESOLUTION



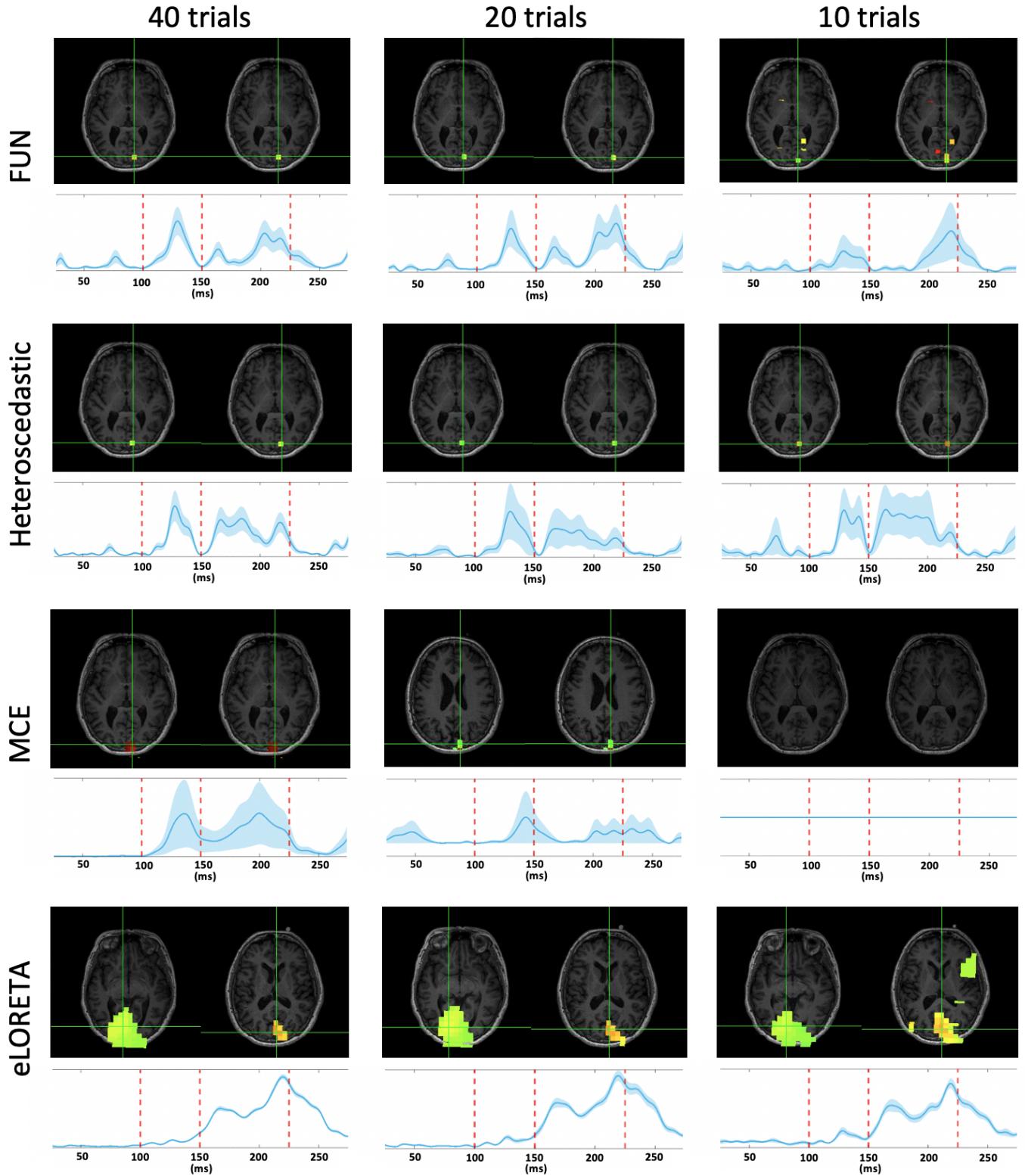
**Fig. 9:** Source reconstruction performance (mean  $\pm$  SEM) of the three different noise learning schemes for data generated by a realistic lead field matrix. Generated sensor signals were superimposed by either full-structure or heteroscedastic noise covering a wide range of SNRs. Performance was measured in terms of the earth mover's distance (EMD), time-course correlation error, F1-measure and Euclidean distance (EUCL) in (mm) between each simulated source and the reconstructed source with highest maximum absolute correlation.



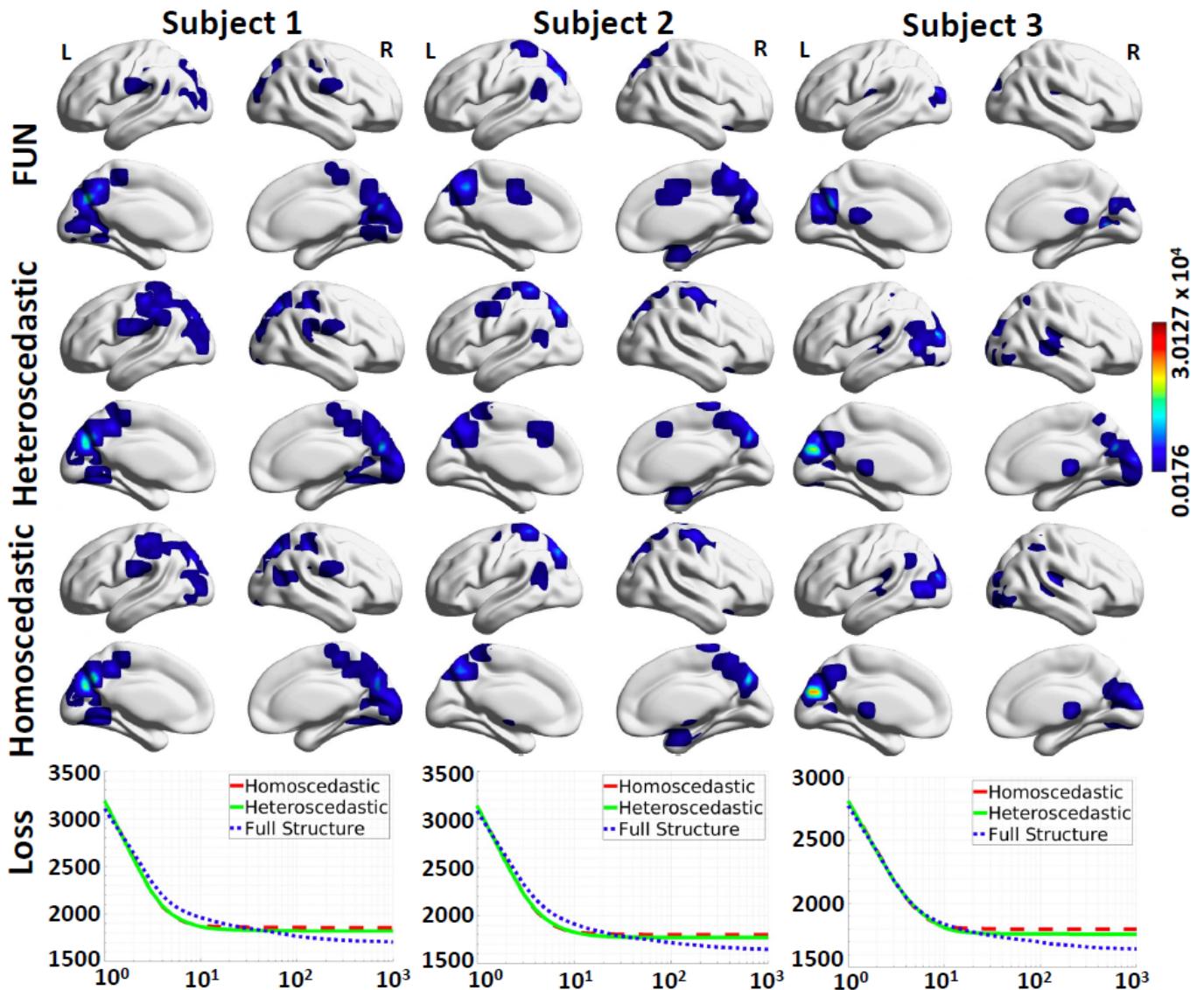
**Fig. 10:** Accuracy of the noise covariance matrix reconstruction incurred by three different noise learning approaches assuming homoscedastic (red), heteroscedastic (green) and full-structural (blue) noise covariances. The ground-truth noise covariance matrix is either full-structure (upper row) or heteroscedastic diagonal (lower row). Performance is assessed in terms of the Pearson correlation between the entries of the original and reconstructed noise covariance matrices,  $\Lambda$  and  $\hat{\Lambda}$ , denoted by  $\Lambda^{\text{sim}}$  (left column). Shown is the similarity error  $1 - \Lambda^{\text{sim}}$ . Further, the normalized mean squared error (NMSE) between  $\Lambda$  and  $\hat{\Lambda}$ , defined as  $\text{NMSE} = \|\hat{\Lambda} - \Lambda\|_F^2 / \|\Lambda\|_F^2$  is reported (right column).



**Fig. 11:** Auditory evoked field (AEF) localization results from one representative subject for different numbers of trial averages using FUN learning, heteroscedastic Champagne, MCE and eLORETA. All reconstructions of FUN learning algorithm show focal sources at the expected locations of the auditory cortex. Even when limiting the number of trials to as few as 10 reconstruction result of FUN learning are accurate, while it severely affects the reconstruction performance of competing benchmark methods.



**Fig. 12:** Localization and time series results of visual evoked field (VEF) activity for a single subject using FUN and benchmarks. Comparing with MCE and eLORETA, FUN shows accurate localization capability. Furthermore, FUN detects sharper 2<sup>nd</sup> peaks when compared to the heteroscedastic noise-learning Champagne, which is consistent with the sharp response of the VEF. The results obtained by FUN are robust across different SNRs/numbers of trial averages. For additional benchmark results, please see supplementary Fig. 8.



**Fig. 13:** Localization of resting-state brain activity for three subjects using FUN and the heteroscedastic and homoscedastic noise learning variants of Champagne. The source variance patterns estimated by each algorithm are projected onto the cortical surface. The convergence behaviour of all three noise estimation approaches is also shown in terms of the negative log-likelihood cost function. FUN converges to better minima when compared to these benchmarks.