

# FUNCTIONAL RANDOM EFFECTS MODELING OF BRAIN SHAPE AND CONNECTIVITY

BY EARDI LILA<sup>1</sup> AND JOHN A.D. ASTON<sup>2</sup>

<sup>1</sup>*Department of Biostatistics, University of Washington*

<sup>2</sup>*Statistical Laboratory, University of Cambridge*

We present a statistical framework that jointly models brain shape and functional connectivity, which are two complex aspects of the brain that have been classically studied independently. We adopt a Riemannian modeling approach to account for the non-Euclidean geometry of the space of shapes and the space of connectivity that constrains trajectories of co-variation to be valid statistical estimates. In order to disentangle genetic sources of variability from those driven by unique environmental factors, we embed a functional random effects model in the Riemannian framework. We apply the proposed model to the Human Connectome Project dataset to explore spontaneous co-variation between brain shape and connectivity in young healthy individuals.

**1. Introduction.** Human brains differ in their structural and functional organization (Gilmore, Knickmeyer and Gao, 2018). While there is a long history of trying to relate either structural or functional brain features to human aspects, such as behavioral and cognitive variables (for recent examples, see, e.g., Xia et al., 2018; Zhang et al., 2019), more recently, increasing attention has been drawn to the problem of understanding how brain structure and function are related to each other (Bullmore and Sporns, 2009).

In this work, we introduce a statistical framework that allows us to estimate patterns of co-variation between brain structure and function, while disentangling co-variation due to genetic and environmental factors. We describe the brain structural organization of an individual with a surface encoding the *brain shape*, that is the geometry of the highly convoluted outermost layer of the brain, called the cerebral cortex. We describe the brain functional organization of an individual with a network that has spatial nodes located on the cerebral cortex. The strength of the network edges is estimated by a measure of pairwise statistical dependence (e.g., correlation) between the neuronal activity associated with the network nodes. The estimated network is a representation of the subject's *brain functional connectivity*. Figure 1 provides an illustration of this setting.

We apply the proposed methodology to 1003 young adults in the Human Connectome Project (HCP) dataset (Glasser et al., 2013) with the aim of exploring spontaneous modes of genetically-driven and environmentally-driven co-variation in the brain structure and function of healthy individuals.

From a methodological perspective, our work can be contextualized within the Object Data Analysis framework (Marron and Alonso, 2014) as both shapes and connectivity networks are complex data objects, living on functional non-Euclidean spaces, where classical Functional Data Analysis approaches (Ramsay and Silverman, 2005) fail to preserve the geometry of these spaces. In order to enforce physiologically valid shape trajectory estimates, we represent brain shapes through diffeomorphic deformation functions of the ambient space. We represent brain connectivity by means of covariance functions, which must be non-negative

---

*Keywords and phrases:* Functional Data Analysis, Variance component models, Mixed effects models, Neuroimaging.

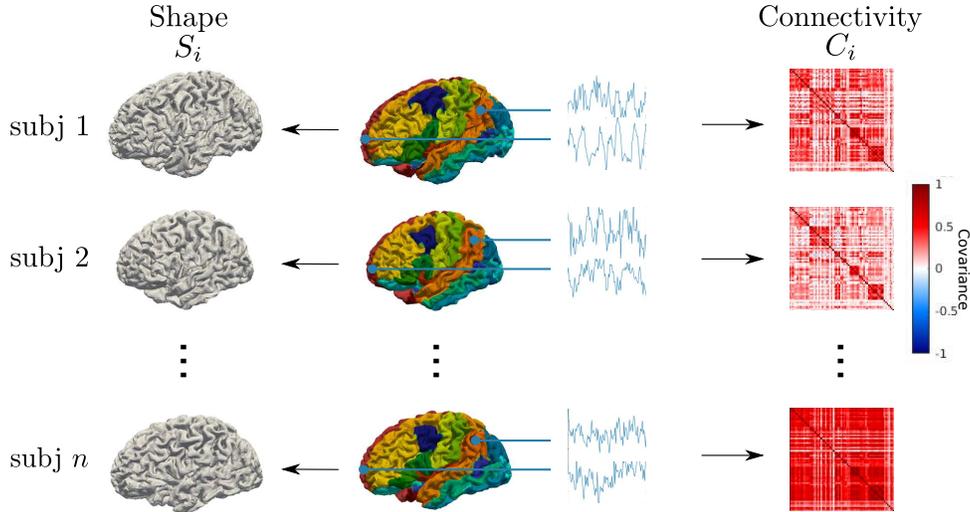


Fig 1: In the central panel, we show the subject-specific surfaces encoding the geometry of the cerebral cortex, as reconstructed from the MRI scans. Moreover, we can see the fMRI time-series, describing the neuronal activity of a dense set of 64K points on the cerebral cortex. The color map on the brain surfaces describes a parcellation atlas, which defines 68 regions of the brain in correspondence across subjects. Within each region, an average time-series is computed. These are then used to compute the  $68 \times 68$  covariance matrices on the right panel, describing the functional connectivity of each subject. On the left panel, a representation of only the shape of the brain surfaces. The shapes on the left panel and the covariances on the right panel are the object-data of our statistical analysis.

definite. Our approach tackles diffeomorphic constraints and non-negative definiteness constraints in a Riemannian framework, i.e., by tangent space mapping through Riemannian logarithmic maps.

Within the proposed Riemannian modeling framework, we define a multi-variable trait *variance component model* that exploits the relatedness structure among individuals to disentangle co-variation in shape and connectivity that is due to genetic sources and environmental sources. The proposed model can be regarded as an extension of the classical single-trait and bivariate-trait variance component models in Amos (1994); Almasy, Dyer and Blangero (1997) and is formulated as a multivariate mixed effects model on the Karhunen–Loève basis coefficients of the tangent space coordinates. Related models have already been employed in the neuroimaging literature to estimate the impact of genes and environment on structural brain development (Lenroot and Giedd, 2008) or to model subject-specific heterogeneity in functional connectivity (Fiecas et al., 2017). However, such studies tend to focus on either structural or functional features.

Mixed effects models have been successfully extended to the setting of functional data in a linear space, to account for non-parametric fixed and random effects (see, e.g., Shi, Weiss and Taylor, 1996; Guo, 2002; Wu and Zhang, 2002; Qin, 2005; Morris and Carroll, 2006; Chen and Wang, 2008; Zhou, Huang and Carroll, 2008; Reimherr and Nicolae, 2016; Liu, Wang and Cao, 2017; Scheipl, Staicu and Greven, 2015). Functional models that incorporate genetic information, without explicitly relying on a mixed effects model have also been formulated. For instance, Kirkpatrick and Heckman (1989) introduce a model to separate genetic functional traits and Luo et al. (2019) propose a model that is able to dissect genetic and environmental effects of functional data in a twin study design. The proposed model applies

to linear functional data and is formulated as a functional structural equation model. An extension to functional data over two-dimensional domains and living in a linear function space, such as cortical thickness data mapped onto a spherical domain, has been proposed in Risk and Zhu (2019). In the non-Euclidean framework of our analysis, the variance component model approach allows for more flexible relatedness structure among individuals, possibly estimated from Single Nucleotide Polymorphism (SNP) data (see, e.g., Dahl et al., 2016).

The HCP dataset, which motivates this work, includes Magnetic Resonance Imaging (MRI) and resting-state functional MRI (fMRI) scans. The MRI scans are used to reconstruct surface models of the cerebral cortex geometry. The time-variant fMRI signals are used to estimate a spatial covariance structure on the cerebral cortex, describing how the different parts of the cerebral cortex co-activate in time, namely an estimate of the functional connectivity. An illustration of the MRI and fMRI components of the data is provided in Figure 1. Moreover, the pedigree of the HCP cohort is available and includes monozygotic twins, dizygotic twins, full siblings, half-siblings, and unrelated individuals. This family structure allows the variance component model to disentangle genetic and environmental co-variation between brain shape and connectivity.

*Statistical analysis of shapes and covariances.* In the neuroimaging literature, brain shape is usually modeled using a few descriptors of shape, such as cortical volume or the area of pre-defined sets of brain regions (Im et al., 2008; Hazlett et al., 2017). In the statistical literature, a non-exhaustive list of shape analysis methodologies based on discrete representations includes landmark-based shape representations (see, e.g., Dryden and Mardia, 2016), skeletal shape representations (Pizer et al., 2013), dihedral angles representations (Eltzner, Huckemann and Mardia, 2018) and projective shape spaces (Mardia and Patrangenaru, 2005). In the continuous setting of curves and surfaces, global parametrizing functions have been adopted to represent these objects.

Representing curves and surfaces with their parametrizing functions, equipped with an  $L^2$  norm, leads to unnatural trajectories in the space of shapes. Instead, a successful approach consists of equipping the space of parametrizing functions with an Elastic Riemannian metric and defining parametrization-invariant representations. The resulting space and associated metric lead to naturally looking geodesic trajectories in the space of curves (Kurtek et al., 2012; Su et al., 2014) and surfaces (Kurtek et al., 2011; Jermyn et al., 2012, 2017).

Of particular importance to this work is a class of representation models for surfaces that do not require the computation of parametrizing functions. These represent surfaces with diffeomorphic deformation functions of the ambient space  $\mathbb{R}^3$  (see, e.g., Vaillant et al., 2004; Charon and Trounev, 2014; Arguillère, Miller and Younes, 2016; Younes, 2010). Such an approach is well suited to the neuroimaging setting because constraining the deformation functions to be diffeomorphic results in a shape space that contains anatomically plausible shapes and excludes, for instance, self-intersecting surfaces. Statistical analysis can then be performed on the non-linear manifold of diffeomorphic functions by exploiting a tangent space expansion to find a linear representation of the data.

In a similar spirit to shape analysis, the statistical analysis of samples that are covariances also involves a non-Euclidean type analysis. In fact, the neuroimaging community has often approached the problem by performing multivariate analysis on vectorizations of the covariances (Smith et al., 2015; Xia et al., 2018). However, such an approach fails to guarantee that linear extrapolations of the data belong to the space of covariances, i.e., that they are positive semi-definite objects. In other words, a signal with the estimated extrapolated ‘covariance’ may not exist. Such an issue can be overcome by defining an appropriate Riemannian metric on the space of covariances.

To this purpose, different metrics have been proposed. For instance, Pennec, Fillard and Ayache (2006) introduce an affine invariant Riemannian metric, while Arsigny et al. (2006)

introduce a log-Euclidean metric based on the matrix exponential and matrix logarithm functions. Dryden, Koloydenko and Zhou (2009) introduce a metric that can deal with rank deficient covariance matrices, and its extension to covariance operators has been proposed in Pigoli et al. (2014). As shown in Dryden, Koloydenko and Zhou (2009), different metrics lead to different geodesic trajectories in the space of covariances. While these trajectories are easy to visualize for lower dimensional covariances, in our high-dimensional setting, these differences are more difficult to appreciate. Therefore, our choice of the metric is mostly driven by computational efficiency arguments, and in particular by closed-form solutions of the geodesic mean. Statistical analysis can then be performed on tangent space projections of the covariances, computed through the Riemannian logarithmic map, which offer a convenient linear parametrization. We then use the tangent coordinates in the space of shapes and the space of covariances to find maximally associated modes of variation in shape and connectivity while decomposing the genetic and unique environmental variance contributions.

The rest of the paper is organized as follows. In Section 2, we introduce the Riemannian modeling framework and the variance component model. In Section 3 we present the implementation details of the proposed model. We then apply the proposed model to the HCP dataset and present the results in Section 4. We finally give some concluding remarks in Section 5. The simulations validating the variance component model are postponed to the appendix.

**2. Mathematical description of the model.** Consider a sample of  $n$  pairs of observations  $\{(S_i, C_i) : i = 1, \dots, n\}$ . Here,  $(S_i \subset \mathbb{R}^3)$  are two-dimensional surfaces, embedded in  $\mathbb{R}^3$ , representing the cerebral cortex geometries. The functions  $(C_i : S_i \times S_i \rightarrow \mathbb{R})$  are covariance functions representing the associated functional connectivity on the subject’s cerebral cortex. Moreover, we assume statistical relatedness between the  $n$  samples; in our application this is due to family-based genetic associations. The aim of this section is to introduce the proposed statistical framework for the analysis of the objects  $\{(S_i, C_i) : i = 1, \dots, n\}$ .

As previously mentioned, both the space of brain geometries and that of brain connectivity are non-Euclidean spaces, introducing additional challenges in the definition and estimation of the co-variation structure. In Section 2.1, we first give a brief conceptual description of the Riemannian approach to modeling shape and connectivity spaces, and then follow by introducing the variance component model. We detail our choices of the representation models and metrics, for the shape and connectivity spaces, in Section 2.2 and Section 2.3, respectively.

*2.1. Functional random effects modeling of shape and connectivity.* Let  $S_0 \subset \mathbb{R}^3$  be a template surface. We assume there exists a one-to-one correspondence between each of the points on  $S_0$  and those on a surface  $S_i$ . The role of the template is two-fold here. The template is a surface representing an “average” geometric shape of the population that allows us to model shapes as functions that are  $\mathbb{R}^3$  deformations of the reference template. Moreover, the template plays the role of a common reference domain where the subject-specific covariance functions can be mapped.

For a fixed template, we represent each surface  $S_i$  with an associated deformation function  $\gamma_i : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  such that  $\gamma_i(S_0) = S_i$ . These deformations are diffeomorphic functions, i.e., they are smooth one-to-one functions with smooth inverse. The space of deformations is formally equipped with a Riemannian metric and the diffeomorphic functions are projected onto the tangent space centered at the identity function. A set of tangent space coordinates  $v_1^S, \dots, v_n^S$  is then used to represent the surfaces  $S_1, \dots, S_n$ . We assume that each function  $v_i^S$  can be expressed in terms of a common basis expansion

$$v_i^S = \sum_{j=1}^{\infty} A_{i,j}^S \psi_j^S,$$

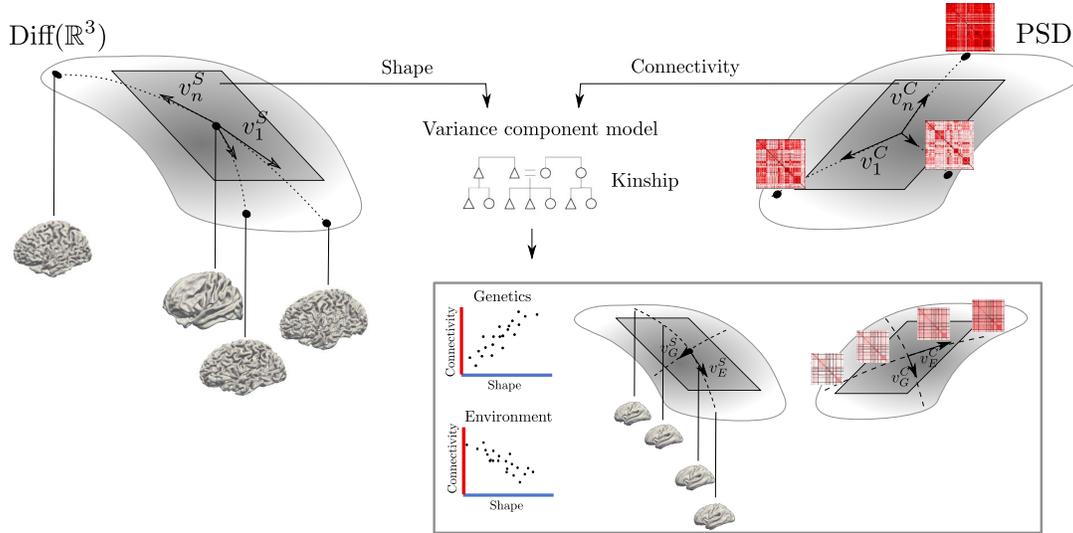


Fig 2: This is an illustration of the proposed statistical analysis framework. We represent shapes with diffeomorphic deformations of the ambient space. The depiction of the space of diffeomorphic functions and that of covariances highlights their non-Euclidean structure. We rely on a tangent space expansion to derive linear parametrizations of these non-Euclidean spaces. The linear tangent coordinates of shape ( $v_i^S$ ) and connectivity ( $v_i^C$ ) are then jointly used to define a variance component model that exploits the kinship structure among the samples to estimate pairs of tangent coordinates that are highly correlated due to genetic factors ( $v_G^S, v_G^C$ ) or environmental factors ( $v_E^S, v_E^C$ ). The shape and covariance non-linear trajectories associated with the estimated pairs of tangent coordinates can finally be computed to display the results in terms of elements of the shape and connectivity spaces.

where  $\psi_j^S$  is the  $j$ th basis element and  $A_{i,j}^S$  is the coefficient of the  $i$ th sample associated with the  $j$ th basis.

The covariance functions ( $C_i$ ) also belong to a non-Euclidean space, which is the cone of positive semi-definite covariance functions. Moreover, they are defined on sample-specific spatial domains ( $S_i$ ). The previously defined deformation functions ( $\gamma_i$ ) can be used to map a covariance  $C_i$  onto the template  $S_0$ , defining  $C_i^0(x, y) := C_i(\gamma_i^{-1}(x), \gamma_i^{-1}(y))$ , with  $x, y \in S_0$ . This leads to a set of ‘spatially normalized’ covariance functions  $C_i^0 : S_0 \times S_0 \rightarrow \mathbb{R}$ . As detailed in Section 2.3, prior to the definition of the Riemannian metric we reduce the covariance functions into finite-dimensional positive-definite covariance matrices. We omit the details here to keep the notation simple, however, it should be noted that this step has implications on the class of metrics that we will be able to adopt (see, e.g., Pigoli et al., 2014, for a discussion). The data ( $C_i^0$ ) are then projected onto the tangent space centered at the geodesic mean and the associated tangent space coordinates  $v_1^C, \dots, v_n^C$  are used to represent  $C_1^0, \dots, C_n^0$ . As for the shape tangent coordinates, we assume ( $v_i^C$ ) can be expressed in terms of a common basis expansion

$$v_i^C = \sum_j A_{i,j}^C \psi_j^C,$$

with  $\psi_j^C$  the  $j$ th basis element and  $A_{i,j}^C$  the coefficient of the  $i$ th sample associated with the  $j$ th basis.

A simple approach to controlling for a set of known confounding variables ( $z_i \in \mathbb{R}^l$ ) consists of modeling the coefficients  $\mathbb{E}[A_{i,j}^S | z_i]$  and  $\mathbb{E}[A_{i,j}^C | z_i]$  through a regression analy-

sis. The conditional expected values  $\mathbb{E}[A_{i,j}^S|z_i]$  and  $\mathbb{E}[A_{i,j}^C|z_i]$  are related to  $\mathbb{E}[v_i^S|z_i]$  and  $\mathbb{E}[v_i^C|z_i]$  by the following equations:

$$\mathbb{E}[v_i^S|z_i] = \sum_j \mathbb{E}[A_{i,j}^S|z_i] \psi_j^S, \quad \mathbb{E}[v_i^C|z_i] = \sum_j \mathbb{E}[A_{i,j}^C|z_i] \psi_j^C$$

The estimated effects of the confounders can then be removed from the tangent space coordinates.

In practice, we choose appropriate truncation levels  $p_S$  and  $p_C$ , and rely on the finite-dimensional approximations

$$v_i^S \approx \sum_{j=1}^{p_S} A_{i,j}^S \psi_j^S, \quad v_i^C \approx \sum_{j=1}^{p_C} A_{i,j}^C \psi_j^C.$$

While any basis in the space of the tangent coordinates is a valid choice, in our application we adopt a principal component basis, due to its well known best linear approximation property.

Before introducing our variance component model, we briefly recall the matrix normal distribution MVN, which generalizes the multivariate normal distribution to matrix-valued random variables. A  $n \times p$  random matrix  $R$  has a matrix normal distribution  $\text{MVN}(M, U, V)$  if and only if  $\text{vec}(R)$  has a multivariate normal  $\mathcal{N}(\text{vec}(M), V \otimes U)$ , where  $\text{vec}$  is the column-wise vectorization operator, and  $\otimes$  denotes the Kronecker product. The matrix normal distribution is characterized by three parameters that are a  $n \times p$  mean matrix  $M$ , a  $n \times n$  matrix  $U$  (modeling the covariance structure between the rows of the random matrix) and a  $p \times p$  matrix  $V$  (modeling the covariance structure between the columns of the random matrix).

Consider now the  $n \times p_S$  coefficient matrix  $(A^S)_{ij} = A_{i,j}^S$  and the  $n \times p_C$  coefficient matrix  $(A^C)_{ij} = A_{i,j}^C$ . We propose a joint model for shape and connectivity in terms of their  $n \times (p_C + p_S)$  coefficient matrix  $A$  that contains both the features described in  $A^S$  and  $A^C$ , i.e.,

$$A := [A^S, A^C].$$

If at this stage we were interested in understanding co-variation between brain shape and connectivity, we could, for instance, perform canonical correlation analysis on the scores matrices  $A^S$  and  $A^C$ , or alternatively, perform an angle-based joint and individual variation analysis (Feng et al., 2018; Carmichael et al., 2019). In fact, several related approaches have been proposed to integrate structural and functional neuroimaging data (see, e.g., Franco et al., 2008; Sui et al., 2011; Xue et al., 2015). However, the estimated joint variation components would be an average of the co-variation that is due to genetic and environmental factors. Therefore, our next step is defining a model that separates genetic and environmental variability.

Let  $K_n$  be a  $n \times n$  matrix of relatedness coefficients, that is, a correlation structure, among the  $n$  subjects, encoding genetic relatedness. We assume this is known. In practice, it can be estimated from a pedigree or from genomic data (Lange, 2002; Kang et al., 2010; Wang, Sverdlov and Thompson, 2017). Let  $I_n$  be the  $n \times n$  identity matrix, which as opposed to  $K_n$ , encodes ‘unrelatedness’ between subjects. Let the  $(p_S + p_C) \times (p_S + p_C)$  matrices  $\Sigma_G$  and  $\Sigma_E$  denote respectively the unknown genetic and environmental covariance structure across the columns of  $A$ . We model  $A$  as

$$\begin{aligned} (1) \quad & A | G, E = XB + G + E \\ & G \sim \text{MVN}(0_{n \times (p_S + p_C)}, K_n, \Sigma_G), \\ & E \sim \text{MVN}(0_{n \times (p_S + p_C)}, I_n, \Sigma_E), \end{aligned}$$

where  $G$  and  $E$  are independent,  $X$  denotes a  $n \times s$  fixed design matrix and  $B$  a  $s \times (p_S + p_C)$  matrix of fixed effects (i.e. not random). In our final application, we do not include a fixed effect term and focus only on the variance component terms  $G$  and  $E$ .

The proposed model exploits a known covariance structure  $K_n$ , across the samples, to additively decompose the covariance structure across the traits into two components, that are  $\Sigma_G$  and  $\Sigma_E$ . In our application,  $K_n$  is chosen to reflect the pairwise family relatedness across samples. Therefore, we refer to  $\Sigma_G$  as the covariance component of the traits that is due to additive genetic contributions, while we refer to  $\Sigma_E$  as the covariance structure that is due to unique environmental contributions. In practice,  $\Sigma_G$  and  $\Sigma_E$  are estimated with a Restricted Maximum Likelihood (REML) approach, as detailed in Section 3.

The model proposed is a multi-variable trait extension of the polygenic quantitative trait variance component model (Amos, 1994; Almasy, Dyer and Blangero, 1997). This is also related to multivariate linear mixed models applied in genome-wide association studies (Zhou and Stephens, 2014; Dahl et al., 2016). In our model, the multivariate traits are tangent space descriptors of brain shape and connectivity.

The matrices  $\Sigma_G$  and  $\Sigma_E$  can be written as

$$\Sigma_G = \begin{bmatrix} \Sigma_G^{S,S} & \Sigma_G^{S,C} \\ \Sigma_G^{C,S} & \Sigma_G^{C,C} \end{bmatrix}, \quad \Sigma_E = \begin{bmatrix} \Sigma_E^{S,S} & \Sigma_E^{S,C} \\ \Sigma_E^{C,S} & \Sigma_E^{C,C} \end{bmatrix},$$

with  $\Sigma_G^{S,S}$  a  $p_S \times p_S$  matrix,  $\Sigma_G^{C,C}$  a  $p_C \times p_C$  matrix,  $\Sigma_G^{S,C}$  a  $p_S \times p_C$  matrix, and  $\Sigma_G^{C,S}$  a  $p_C \times p_S$ . The environmental components are defined similarly. These matrices represent the covariances between and within low-dimensional representations of the tangent space coordinates of shape and connectivity. In this form, they are not themselves very informative. Therefore, we propose to look at maximally correlated genetic and environmental modes of co-variation between the geometry of the cerebral cortex and associated connectivity by computing

$$(\theta_G^S, \theta_G^C) = \left\{ \arg \max_{\theta, \eta} \theta^T \Sigma_G^{S,C} \eta : \theta^T \Sigma_G^{S,S} \theta = 1, \eta^T \Sigma_G^{C,C} \eta = 1 \right\}$$

and

$$(\theta_E^S, \theta_E^C) = \left\{ \arg \max_{\theta, \eta} \theta^T \Sigma_E^{S,C} \eta : \theta^T \Sigma_E^{S,S} \theta = 1, \eta^T \Sigma_E^{C,C} \eta = 1 \right\},$$

where  $\theta_G^S, \theta_E^S$  are  $p_S$ -dimensional vectors, while  $\theta_G^C, \theta_E^C$  are  $p_C$ -dimensional vectors. Constraints of the type  $\theta^T \theta = 1, \eta^T \eta = 1$  are also valid choices. The pairs  $(\theta_G^S, \theta_G^C)$  and  $(\theta_E^S, \theta_E^C)$  represent the maximally correlated modes of co-variation in shape and connectivity that are due to genetic and non-genetic factors. Subsequent modes of co-variation can also be estimated by maximizing the same objective functions while imposing orthogonality constraints with respect to the previously computed components.

We then compute the tangent space coordinates associated with these modes of co-variation, i.e.,

$$(v_G^S, v_G^C) := \left( \sum_{j=1}^{p_S} \theta_{j,G}^S \psi_j^S, \sum_{j=1}^{p_C} \theta_{j,G}^C \psi_j^C \right), \quad (v_E^S, v_E^C) := \left( \sum_{j=1}^{p_S} \theta_{j,E}^S \psi_j^S, \sum_{j=1}^{p_C} \theta_{j,E}^C \psi_j^C \right).$$

Finally, the co-variation structure between brain shape and connectivity, due to genetic factors, can be visualized by computing the surfaces and covariance functions, in their respective curved spaces, that are the Riemannian exponentials of the pairs of elements  $(-c_1 v_G^S, -c_2 v_G^C)$  and  $(c_1 v_G^S, c_2 v_G^C)$ , with  $c_1$  and  $c_2$  appropriately chosen positive constants. Analogously, we

can visualize the environmental co-variation structure between brain shape and connectivity, by computing the surfaces and covariance functions that are the Riemannian exponentials of the pairs of elements  $(-c_1 v_E^S, -c_2 v_E^C)$  and  $(c_1 v_E^S, c_2 v_E^C)$ . The overall analysis is depicted in Figure 2.

*2.2. Brain shape modeling.* We model shapes as deformations of the template  $S_0$  and we define a Riemannian metric on the space of such deformations. In our application setting, this approach allows us to constrain our statistical estimates to be anatomically plausible shapes.

In detail, we introduce a diffeomorphic operator  $\varphi$  mapping a sufficiently smooth Hilbert space  $(\mathcal{V}, \|\cdot\|_{\mathcal{V}})$  onto a group  $\mathcal{G}$  of diffeomorphic functions. Formally, the space  $\mathcal{V}$  is the tangent space of  $\mathcal{G}$  at the identity map, and  $\varphi$  is the associated exponential map. We define  $\mathcal{G}$  as follows. Let  $\{v_t \in \mathcal{V} : t \in [0, 1]\}$  be a time-dependent  $\mathcal{V}$ -valued process such that  $\int_0^1 \|v_t\|_{\mathcal{V}}^2 dt < \infty$ . Then, the solution  $\phi_v : [0, 1] \times \mathbb{R}^3 \rightarrow \mathbb{R}^3$ , at time  $t = 1$ , of the Ordinary Differential Equation (ODE)

$$(2) \quad \frac{\partial \phi_v}{\partial t}(t, x) = v_t \circ \phi_v(t, x) \quad t \in [0, 1], x \in \mathbb{R}^3,$$

with initial condition  $\phi_v(0, x) = x$ , is a smooth diffeomorphic deformation of  $\mathbb{R}^3$  (see, e.g., Younes, 2010). The group  $\mathcal{G}$  consists of all such solutions of equation (2).

Equation (2) allows us to parameterize a diffeomorphic deformation  $\phi_v(1, \cdot)$  (and therefore a shape  $\phi_v(1, S_0)$  that preserves the topology of  $S_0$ ) with a time-variant vector-field  $\{v_t \in \mathcal{V} : t \in [0, 1]\}$ . We then model the space of the time-variant vector-fields by defining  $\{v_t : t \in [0, 1]\}$  to be a time-variant vector field which minimizes the quantity  $\int_0^1 \|v_t\|_{\mathcal{V}}^2 dt$ , for a given initial vector field  $v_0$  (Miller, Trouvé and Younes, 2006). Finally, the diffeomorphic operator is defined to be  $\varphi_{v_0}(x) := \phi_v(1, x)$ , where  $v_0 \in \mathcal{V}$  is the initial vector field generating  $\{v_t : t \in [0, 1]\}$ , and  $\phi_v$  is the solution of the ODE (2) for the computed  $\{v_t : t \in [0, 1]\}$ .

With the notation of the previous section, we define  $\gamma_i := \varphi_{v_i^S}$  with  $i = 1, \dots, n$ . The initial vector fields  $(v_i^S)$  are estimated by minimizing a penalized mismatching functional, the details of which are left to Section 3. What is important from a statistical perspective is that a surface  $S_i$ , which belongs to a curved space, can now be represented by an element  $v_i^S$  of the linear function space  $\mathcal{V}$ , where  $v_i^S$  is such that  $S_i = \varphi_{v_i^S}(S_0)$ .

*2.3. Brain connectivity modeling.* The spatially normalized covariance functions  $(C_i^0)$  associated with densely observed functional data on two-dimensional domains have the additional issue of being voluminous. Therefore, we first expand them into a finite functional basis  $\{b_j\} \subset L^2(S_0)$ , i.e.,

$$C_i^0(x, y) \approx \sum_{j=1}^K \sum_{l=1}^K C_{ijl}^K b_j(x) b_l(y),$$

where  $(C_{ijl}^K)$  is the  $K \times K$  covariance matrix that is a reduced representation of the covariance function  $C_i^0$ .

The functions  $(b_j)$  are estimated from the data. A popular choice in neuroimaging is a set of indicator functions that are constant within connected regions of the cortical surface. This set of connected regions, also known as brain parcellation, defines functional sub-units on the cortical surface. An example of a parcellation is given by the color maps in the central panel of Figure 1. An alternative popular approach consists of estimating such a basis from an Independent Component Analysis of the fMRI data (Calhoun, Liu and Adalı, 2009). Here we adopt the former approach.

Equipping the space of positive definite  $K \times K$  covariance matrices with the  $L^2$  distance results in variations around the mean that may not belong to the space of positive definite objects. We therefore define a Riemannian metric on the space of symmetric positive-definite matrices by defining a smoothly varying scalar product on the tangent space, which is the linear space of  $K \times K$  symmetric matrices. Such a Riemannian metric defines a geodesic distance on the space of covariance matrices that is given by the length of the shortest curve connecting any two covariances. Pennec, Fillard and Ayache (2006) introduce an affine invariant Riemannian metric for positive-definite matrices that induces the distance  $d_{\text{Riem}}(C_1, C_2) = \|\log(C_1^{-1/2} C_2 C_1^{-1/2})\|_F$ , where  $C^{-1/2} = V D^{-1/2} V^T$ , with  $C = V D V^T$  denoting its spectral decomposition. A further option is the Cholesky distance  $d_{\text{chol}}(C_1, C_2) = \|\text{chol}(C_1) - \text{chol}(C_2)\|_F$ , where  $L = \text{chol}(C)$  denotes the Cholesky decomposition of a positive-definite matrix  $C = L L^T$ . In Arsigny et al. (2006), the log-Euclidean distance of two positive-definite matrices  $C_1$  and  $C_2$  is defined as  $d_{\log}(C_1, C_2) = \|\log(S_1) - \log(S_2)\|_F$ , where  $\|\cdot\|_F$  is the Frobenius norm and  $\log(\cdot)$  the matrix logarithm, i.e.,  $\log(C) = V \log(D) V^T$ , with  $\log(D)$  denoting the diagonal matrix whose entries are the logarithms of the entries of  $D$ .

We model covariance matrices by equipping the space of covariances with the log-Euclidean metric defined in Arsigny et al. (2006). For such a choice of the Riemannian metric, the geodesic distance is given by

$$d(C_1, C_2) = \|\log(C_1) - \log(C_2)\|_F.$$

The Fréchet mean  $F \in \mathbb{R}^{K \times K}$  of  $(C_i^K)$  is then defined as  $F = \arg \inf_M \sum_{i=1}^n d(C_i^K, M)^2$ . Given  $C$ , a  $K \times K$  symmetric positive-definite matrix, and  $L$  a  $K \times K$  symmetric matrix, the Riemannian exponential and logarithmic maps have the form

$$\text{Exp}_F(L) = \exp(\log(F) + \partial_L \log(F)), \quad \text{Log}_F(C) = \partial_{\log(C) - \log(F)} \exp(\log(F)),$$

where  $\partial_L \log(F)$  and  $\partial_V \exp(L)$  are respectively the differential of the matrix logarithm and the differential of the matrix exponential (Arsigny et al., 2006). The tangent representation  $v_i^C$  of  $C_i^K$  is then given by the  $K \times K$  symmetric matrix  $\text{Log}_F(C_i^K)$ . We also explored the application of other metrics with closed-form solution for the geodesic mean, such as the Cholesky metric and the square-root metric (Dryden, Koloydenko and Zhou, 2009). In our final application, for the chosen covariance size  $K$ , the different metrics did not seem to make a noticeable difference to the subsequent analysis. As mentioned previously, when  $K$  is very large and potentially tends to infinity, the choice of metric might need to reflect this.

### 3. Estimation.

*Shape representation.* The surfaces  $(S_i)$  need first to be registered; i.e., one-to-one correspondence needs to be established between subject-specific dense sets of landmarks  $(x_i^{(i)}) \subset S_i$ . In practice, the landmarks are the vertices of the triangulated surfaces approximating the idealized surface  $S_i$ . In order to avoid burdening the notation, we do not distinguish between idealized surfaces and associated triangulated surfaces, as this will be clear from the context.

The problem of image registration is common to any population analysis of images, however, the choice of the registration model generally depends on the specific application (Zitová and Flusser, 2003). In our application setting, this step is performed by maximizing a measure of structural/functional ‘coherence’ across subjects, while minimizing the amount of distortion introduced by the registration (Fischl, Sereno and Dale, 1999; Yeo et al., 2010; Robinson et al., 2014, 2018).

Given the  $n$  sets of registered landmarks, a template  $S_0$  — represented by the vertices  $\{x_l^{(0)}\} \subset S_0$  — is estimated by means of a Procrustes analysis. Such an analysis allows us to estimate the template while removing translation, size, and rigid rotations from the surfaces  $(S_i)$ . Then, the shape tangent space coordinates  $(v_i^S)$ , associated with the surfaces  $(S_i)$ , are computed by solving a minimization problem of the form

$$(3) \quad v_i^S = \arg \min_{v_i \in \mathcal{V}} \sum_l \|\varphi_{v_i}(x_l^{(0)}) - x_l^{(i)}\|_{\mathbb{R}^3}^2 + \lambda \|v_i\|_{\mathcal{V}}^2, \quad i = 1, \dots, n,$$

where the least-square term measures the similarity of the deformed template  $\varphi_{v_i}(S_0)$  with  $S_i$ . The term  $\|v_i\|_{\mathcal{V}}^2$  can be intuitively understood as a regularizing term that measures the ‘energy’ associated with the deformation. The constant  $\lambda$  controls the trade-off between the empirical and regularizing term.

To obtain an unbiased estimate of the template  $S_0$ , with respect to the defined metric on the group of diffeomorphic functions  $\mathcal{G}$ , we could update the template with  $S_0 \leftarrow \varphi_{\bar{v}}(S_0)$ , where  $\bar{v} := n^{-1} \sum_{i=1}^n v_i^S$ . Subsequently, we could recompute  $\{v_i^S\}$  by solving (3) for the newly estimated template. These steps can then be iterated until convergence. Such a procedure is however prohibitive for computational reasons. Therefore, we fix the template to be the one resulting from the Procrustes analysis.

The space  $\mathcal{V}$  is modeled as a Reproducing Kernel Hilbert Space (RKHS) with a kernel that is a finite sum of Gaussian kernels of the type  $K_\sigma(x, y) = \exp(-\frac{\|x-y\|_{\mathbb{R}^3}^2}{2\sigma^2})I_3$ , for different choices of  $\sigma > 0$ . The minimization problem is approached with a BGFS optimization scheme (Lewis and Overton, 2013). We adopt the implementation in the MATLAB package `fshapetk` (Charlier, Nardi and Trouvé, 2015; Charlier, Charon and Trouvé, 2017).

*Functional connectivity.* Given the  $K \times K$  subject-specific covariance matrices  $(C_i^K)$  (see Figure 1 for an illustration), we first compute the least-square Fréchet mean estimate. Given our choice of the Riemannian metric, this has the closed-form solution (Arsigny et al., 2006)

$$F = \exp \left\{ \frac{1}{n} \sum_{i=1}^n \log C_i^K \right\}.$$

The tangent space coordinates, in a matrix form, are given by  $V_i^C = \text{Log}_F(C_i^K)$ . In practice, to circumvent stability issues related to the numerical computation of the differential of the matrix exponential and logarithm, we perform tangent expansion around the identity, i.e., work with the coordinates

$$V_i^C = \log(C_i) - \log(F).$$

Finally, the tangent space coordinates  $v_i^C$ , defined in Section 2.3, are given by  $v_i^C = \text{vec}_{\text{Sym}}(V_i^C)$ , where  $\text{vec}_{\text{Sym}}(L) = (l_{1,1}, \dots, l_{K,K}, \sqrt{2}l_{1,2}, \dots, \sqrt{2}l_{K-1,K})^T$  is a convenient vectorization operation for the space of symmetric matrices equipped with the Frobenius norm.

*Mixed Effects Model.* We reformulate Model (1), by defining two  $n \times (p_S + p_C)$  independent random matrices

$$U \sim \text{MVN}(0_{n \times (p_S + p_C)}, I_n, I_{p_S + p_C}), \quad V \sim \text{MVN}(0_{n \times (p_S + p_C)}, I_n, I_{p_S + p_C}).$$

Moreover, we define  $K_n^{\frac{1}{2}}$ ,  $\Sigma_G^{\frac{1}{2}}$ , and  $\Sigma_E^{\frac{1}{2}}$  such that  $K_n = K_n^{\frac{1}{2}}(K_n^{\frac{1}{2}})^T$ ,  $\Sigma_G = \Sigma_G^{\frac{1}{2}}(\Sigma_G^{\frac{1}{2}})^T$ , and  $\Sigma_E = \Sigma_E^{\frac{1}{2}}(\Sigma_E^{\frac{1}{2}})^T$ . Then, we can rewrite Model (1) as

$$A|U, V = XB + K_n^{\frac{1}{2}}U(\Sigma_G^{\frac{1}{2}})^T + V(\Sigma_E^{\frac{1}{2}})^T,$$

thanks to the fact that

$$K_n^{\frac{1}{2}} U (\Sigma_G^{\frac{1}{2}})^T \sim \text{MVN}(0_{n \times (p_S + p_C)}, K_n, \Sigma_G), \quad V (\Sigma_E^{\frac{1}{2}})^T \sim \text{MVN}(0_{n \times (p_S + p_C)}, I_n, \Sigma_E).$$

In practice, the kinship matrix  $K_n$  is rank deficient, which is in fact a desirable property as it means that there are highly correlated samples that, intuitively, make it possible to disentangle the genetic and environmental covariance. When two samples are maximally correlated, as for instance in the case of monozygotic twins, the number of rows of  $U$  can be reduced by one, and the two samples can be modeled with the same random effect.

The matrices  $U$  and  $V$  are then vectorized to a multivariate normal vector with independent samples, and the matrices  $K_n^{\frac{1}{2}}$ ,  $(\Sigma_G^{\frac{1}{2}})^T$ , and  $(\Sigma_E^{\frac{1}{2}})^T$  are reshaped accordingly. Finally, the entries of the unknown covariances  $\Sigma_G$  and  $\Sigma_E$  are estimated optimizing the REML criterion with respect to the parametrizing matrices  $\Sigma_G^{\frac{1}{2}}$  and  $\Sigma_E^{\frac{1}{2}}$ . The proposed model has been implemented in R as a wrapper of the function `lmer` in the package `lme4` (Bates et al., 2015). In our application,  $\Sigma_G$  and  $\Sigma_E$  are unstructured covariance matrices. Nevertheless, the proposed model can be easily extended to handle structured covariance matrices  $\Sigma_G$  and  $\Sigma_E$ .

#### 4. Statistical analysis & Results.

*Data & Preprocessing.* This study focuses on the analysis of all 1003 healthy adult subjects, from the S1200 HCP data release (Van Essen et al., 2013), that have complete resting-state fMRI scans. The structural MRI images have been acquired at a resolution of 0.7mm isotropic and the resting-state fMRI images have been acquired at a spatial resolution of 2.0mm isotropic and a temporal resolution of 0.7s. Resting-state fMRI data were acquired in four runs of 15 minutes. During these fMRI scans, the subjects were not performing any explicit tasks. Further details on the acquisition process can be found in Glasser et al. (2013); Smith et al. (2013). An extensive set of subject traits, such as behavioral and demographic covariates, are also provided. Moreover, the HCP dataset includes multiple members of the same families, resulting in a familial relatedness structure across samples.

The MRI and fMRI data have been pre-processed with the minimal pre-processing HCP pipeline (Glasser et al., 2013). In particular, white, pial, and midthickness surfaces of the cerebral cortex are reconstructed. We use the midthickness surfaces, which are surfaces that interpolate the midpoints between the white and pial surfaces, to describe the anatomy of the cerebral cortex and we refer to them as cortical surfaces. The four resting-state fMRI runs have been pre-processed to remove artifactual components in the data (Smith et al., 2013). The fMRI signals associated with neuronal activation on the cerebral cortex have been extracted and mapped onto the cortical surfaces, resulting in the data illustrated in Figure 1.

*Spatial normalization.* The cortical surfaces are given in the form of two closed triangulated surfaces of 32K vertices, describing respectively the geometry of the left and right hemispheres. The 64K vertices have been brought in correspondence across subjects thanks to the application of a multi-modal surface alignment algorithm (Robinson et al., 2014, 2018), which enables surface alignment based on both anatomical and functional features.

In the setting of joint shape and functional modeling, the importance of using the functional component of the data in the alignment procedure has been demonstrated, for instance, in Lila and Aston (2020). In the cited work, the authors propose a functional manifold surface alignment model embedded in the statistical analysis to improve the surface alignment. Here, we rely on the multi-modal spherical alignment of Robinson et al. (2018), where extensive hyper-parameters tuning and validation have already been performed for the dataset in question.

In the next section, statistical shape analysis is performed on the cortical surfaces by treating the surface vertices as anatomical landmarks, given that these are in geometric and functional correspondence across subjects. Note that the defined surfaces cannot yet be regarded as *shapes* (Dryden and Mardia, 2016) since non-physiological features, such as translation and rotations, are still present in the data.

*Shape Analysis.* In order to project the surfaces in the shape-space — which is, to remove Euclidean similarity transformations from the data — we perform Generalized Procrustes Analysis (Dryden and Mardia, 2016) on the anatomical landmarks. This has the effect of removing translation, rigid rotation, and scale from the data, while iteratively estimating an average shape in the shape-space. Translation and rotation are discarded as these components do not have a physiological meaning. Scale, instead, is a positive scalar that does have a physiological meaning and its log-transformation  $l_i$  will be incorporated in the final analysis. We denote with  $S_i$  the brain surfaces projected onto the shape-space and with  $S_0$  the estimated template average shape.

While features of the data that are non-descriptive of shape have now been removed, the landmark description of shapes does not guarantee that linear statistical models of the data, such as PCA, generate topologically valid shape trajectories. In our application, a valid shape trajectory is one that, for instance, does not lead to self-intersecting surfaces (see, e.g., Viallant et al., 2004, for an illustration of the issue).

Therefore, we represent each surface  $S_i$  with a vector field  $v_i^S$  belonging to the linear space  $\mathcal{V}$ . The vector field  $v_i^S$  is computed by solving the minimization problem in (3), leading to a function  $v_i^S$  such that  $S_i \approx \varphi_{v_i^S}(S_0)$ , where  $\varphi$  is the diffeomorphic deformation operator defined in Section 2.2. In practice, two separate vector fields, one for each hemisphere, are estimated independently. We do not denote them separately as such a choice has computational advantages, but no other practical implications.

The RKHS space  $\mathcal{V}$  is defined by a kernel that is the sum of six isotropic Gaussian kernels in  $\mathbb{R}^3$  with standard deviations  $\{8, 4, 2, 1, 0.5, 0.1\}$ . The penalty coefficient in (3) is chosen to be  $\lambda = 10^{-3}$ . These hyper-parameters have been selected by experimenting with a small subset of the full cohort. Computations have been performed on a cluster where each node is equipped with an Intel Xeon E5-2650 2.2GHz 12-core processor with 96GB RAM and four Nvidia P100 GPUs. The minimization algorithm takes approximately 40 minutes for each subject and uses only one GPU. The representing 1003 vector fields are computed in parallel on several nodes, greatly reducing the computation time needed.

We identify a set of demographic confounding variables (height, weight, sex, and age), which are demeaned, and their squares (when the confounder is a continuous variable) are included in the analysis. We then regress the confounders out of the RKHS coefficients representing  $(v_i^S)$  and the log-transformed size coefficients ( $l_i$ ) from the Procrustes Analysis.

We perform functional PCA analysis on  $\{v_i^S\}$  leading to the representation  $v_i^S \approx \bar{v}^S + \sum_{j=1}^{p_S} A_{i,j}^S \psi_j^S$ . In contrast to the idealized basis expansion in Section 2.2, we have a non-zero mean term  $\bar{v}^S = n^{-1} \sum_{i=1}^n v_i^S$ . In fact, due to the prohibitive computational costs incurred in computing  $v_i^S$ , we do not iteratively re-estimate the template until the mean term  $\bar{v}^S$  becomes negligible, as noted in Section 3. We select the truncation level  $p_S = 10$ . This choice is mainly driven by computational limitations in the subsequent analysis. Our final shape representation of the cerebral cortex of the  $i$ th subject is given by  $p_S$  scalar coefficients that are  $A_{i,1}^S, \dots, A_{i,p_S}^S$ , and the log-transformed size coefficient  $l_i$  from the Procrustes analysis.

*Connectivity Analysis.* For each vertex of a surface  $S_i$ , we have four times series (one for each run) describing the resting-state neuronal activation of that location. Thanks to the anatomical and functional correspondence of the surface vertices across subjects, we can

equivalently perform our analysis on the common template surface  $S_0$ . We adopt an atlas of the template  $S_0$  that assigns a label to each vertex of the template, defining a parcellation. For each run, and within each region of the parcellation, we compute a robust spatially averaged time-series that represents the brain activity of that entire region.

Many approaches have been proposed to define parcellation atlases (Fischl et al., 2004; Desikan et al., 2006; Power et al., 2011; Yeo et al., 2011; Van Essen et al., 2012; Wig, Laumann and Petersen, 2014; Gordon et al., 2016; Glasser et al., 2016). See Arslan et al. (2018) for a recent systematic review. Over the years they have tremendously improved in their granularity and ability to incorporate multi-modal imaging to define parcellations of the cortical surface. In this work, we rely on the popular Desikan-Killiany parcellation (Desikan et al., 2006), which defines  $K = 68$  cortical surface regions. More recent parcellations, as the one proposed in Glasser et al. (2016), define up to 360 regions. We opted for a courser parcellation to mitigate the effect of misregistration, which is the effect of small errors in the surface registration step.

For each subject, the  $68 \times 4$  time-series are demeaned and variance normalized. A  $68 \times 68$  covariance matrix is computed for each run and the average covariance across the four runs,  $C_i^K$ , is used to represent the  $i$ th subject functional connectivity. These covariance matrices are sometimes referred to as networks, or connectomes, as they quantify the functional connectivity between network nodes that are regions of the cortical surface.

We perform connectivity analysis in the log-Euclidean framework. As detailed in Section 2.3, we compute the matrix logarithms of  $C_i^K$  and compute the associated coefficients with respect to a Frobenius-orthogonal basis on the space of symmetric matrices. This leads to the computation of a set of representing tangent space coordinates  $\{v_i^C\}$  that are vectors of dimension 2346, which is the number of entries of the upper triangular part of the  $68 \times 68$  covariance matrices.

In addition to those used in the shape analysis step, we identify additional confounding variables that are more directly related to the acquired fMRI signal (Smith et al., 2015) (acquisition reconstruction software version, average subject head motion, systolic/diastolic blood pressure, hemoglobin A1C measured in blood, cube-root of total brain volume, and cube-root of total intracranial volume). The confounding variables are regressed out of the connectivity tangent coordinates  $\{v_i^C\}$  (we do not rename the deconfounded tangent space coordinates).

We perform PCA on the deconfounded tangent space coordinates  $\{v_i^C\}$ , leading to the basis expansion  $v_i^C \approx \bar{v}_i^C + \sum_{j=1}^{p_C} A_{i,j}^C \psi_j^C$ . We truncate the basis expansion at  $p_C = 10$ . As in the shape analysis step, this choice is mainly driven by computational limitations in the subsequent joint shape/connectivity analysis. For the  $i$ th subject, connectivity is finally represented by 10 coefficients that are  $A_{i,1}^C, \dots, A_{i,p_C}^C$ .

*Family relatedness.* The statistical model proposed in Section 2.1 relies on the presence of genetic relatedness among the subjects to disentangle genetic and environmental contributions to the covariation patterns between brain shape and connectivity. Of the 1003 subjects, 1001 had family relatedness information. The dataset consists of unrelated individuals, full siblings, half-siblings, dizygotic twins, and monozygotic twins. Specifically, there are 429 families, with a number of members that range from 1 to 6.

The relatedness matrix  $K_n$ , in the multivariable trait model (1), represents pairwise expected covariance, between subjects, that is due to familial relatedness. We estimate the matrix  $K_n$  as  $K_n = 2\Phi_n$ , with  $\Phi_n$  the matrix of kinship coefficients (Almasy, Dyer and Blangero, 1997; Lange, 2002). We use Solar ([nitrc.org/projects/se\\_linux](http://nitrc.org/projects/se_linux)) to compute the kinship matrix from the HCP family structure data. In larger population studies, an empirical genetic relatedness matrix could be computed from SNP data (see, e.g., Kang et al., 2010; Wang, Sverdlov and Thompson, 2017).

*Joint random-effects modeling.* In the previous steps of the analysis, for the  $i$ th subject, we have derived a vector of scalar variables

$$(4) \quad l_i, A_{i,1}^S, \dots, A_{i,p_S}^S, A_{i,1}^C, \dots, A_{i,p_C}^C,$$

with the first variable being a descriptor of brain size, the subsequent  $p_S$  variables being descriptors of shape, and the final  $p_C$  variables being descriptors of connectivity. The empirical  $(1 + p_S + p_C) \times (1 + p_S + p_C)$  covariance matrix  $\Sigma$ , computed from these variables, describes the pairwise first-order dependencies between such descriptors, hence the co-variation structure between shape and connectivity that is due to both genetic and environmental contributions. We instead want to leverage the familial relatedness matrix  $K_n$  and apply the joint mixed model in Section 2.1 to disentangle the additive covariance components  $\Sigma_G$  and  $\Sigma_E$ , which are respectively due to genetic and environmental contributions.

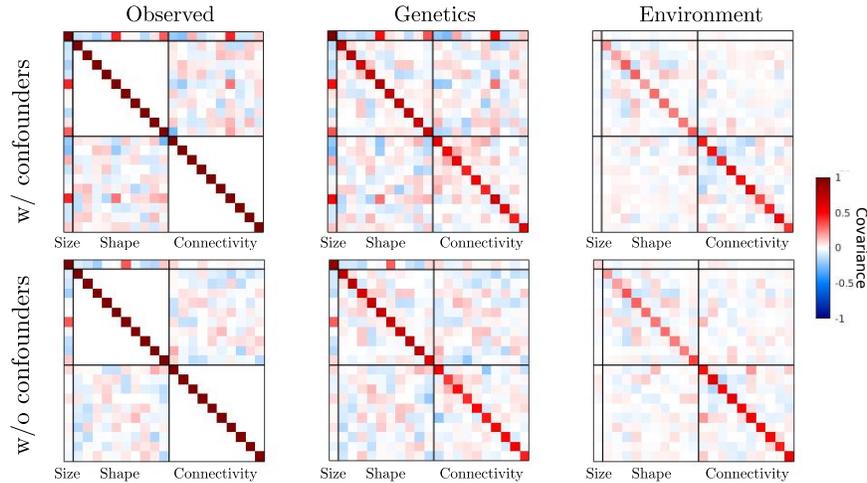


Fig 3: On the left panel, we plot the empirical covariance of the  $n \times (1 + p_S + p_C)$  data matrix of size, shape, and connectivity descriptors. Each descriptor (i.e., each column of the data matrix) has been normalized to have unit standard deviation. On the middle and right panel, we can find the latent covariance components that are due to genetic and environmental factors, as recovered by the variance component model proposed.

The covariances  $\Sigma_G$  and  $\Sigma_E$  are estimated by optimizing the REML criterion associated with Model (1), without a fixed effects term. For our choice of the truncation levels ( $p_S = p_C = 10$ ), the minimization of the REML takes approximately 8 hours. This is the limiting factor in the choice of the truncation levels. Nonetheless, in general, care should be taken in increasing  $p_S$  and  $p_C$ , as the descriptors could start capturing smaller scale variations in shape and connectivity driven by small and unavoidable errors in the surface alignment steps (for an illustration of the issue, see Section S3 in the supplementary material for Lila and Aston, 2020).

*Results.* In Figure 3, we show respectively the estimates of  $\Sigma$  (the empirical covariance structure),  $\Sigma_G$  (the covariance structure due to genetic contributions), and  $\Sigma_E$  (the covariance structure due to unique environmental contributions) of the  $1 + p_S + p_C$  descriptors of size, shape, and connectivity.

We measure the heritability, i.e., the portion of shape and connectivity variability explained by additive genetic contributions as

$$h^2 = \frac{\text{trace}(\Sigma_G)}{\text{trace}(\Sigma_G + \Sigma_E)}.$$

Our overall heritability estimate is  $h^2 = 0.61$ , which means we estimate that 61% of the source of variability in the data (post-PCA) is due to genetic contributions. However, if we measure the heritability of size, shape, and connectivity separately, we obtain the estimates 0.92, 0.71, and 0.47, respectively. These are consistent with recent estimates in the literature (Barber et al., 2021; Pizzagalli et al., 2020), and with the intuition that functional features are more easily affected by life events and environmental factors, and therefore are less heritable. This is also clear from the diagonal entries of  $\Sigma_G$  and  $\Sigma_E$  in Figure 3.

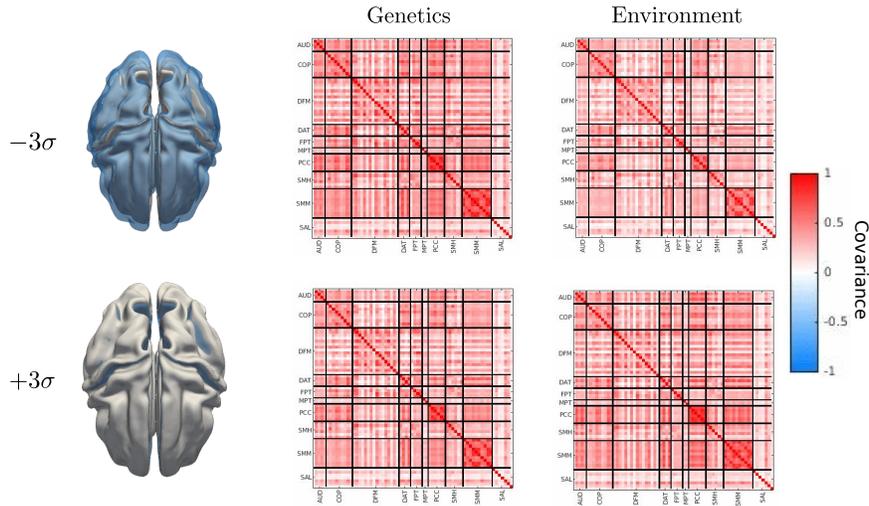


Fig 4: An illustration of the changes in connectivity levels associated with a  $\pm 3\sigma$  variation in the univariate variable log-transformed brain size, where  $\sigma$  is the standard deviation of such a variable. The corresponding anatomical objects are shown as gray-colored surfaces. The blue surfaces represent the fixed average brain and are shown in order to have a fixed reference between plots and be able to appreciate the differences between the two gray surfaces. In the areas where only the gray surface is visible, this is overlaying the blue surface. The associated changes in genetic and environmental connectivity levels, shown in the form of covariance matrices, are minimal. Instead, running the analysis without removing the effect of confounders displayed a general increase in connectivity that is associated with an increase in brain size. In order to improve the readability of the covariance plots, we have clustered their nodes in pre-identified communities (Power et al., 2011): auditory network (AUD), cingulo-opercular network (COP), default mode network (DMN), dorsal attention network (DAT), frontoparietal network (FPT), medial parietal network (MPT), somatosensory/motor network (SMH & SMM), and salience network (SAL).

These heritability estimates are a byproduct of our analysis. The main contribution of our work is instead to quantify the dependence structure between anatomical and functional features and display this dependence in terms of actual variations of the neurobiological object considered. In Figure 4, we show the results of a linear regression model, between

brain size and connectivity, formulated using the estimated genetic and environmental covariances. After the linear model is fitted, we display the variation in connectivity associated with a  $\pm 3\sigma$  variation of the log-transformed size variable, where  $\sigma$  is its standard deviation. The results do not display large associated variations in connectivity. However, when running the analysis without removing the confounders from the size and connectivity descriptors, we observe that a general increase in connectivity is associated with an increase in brain size, both in the genetic and environmental components. This association is driven by the confounding factors. In Figure 5, we display the results of the CCA between shape descriptors and connectivity descriptors. Specifically, we compute the tangent coordinates  $(-3\sigma_G^S v_G^S, -3\sigma_G^C v_G^C)$  and  $(3\sigma_G^S v_G^S, 3\sigma_G^C v_G^C)$ , for the genetic contribution, and  $(-3\sigma_E^S v_E^S, -3\sigma_E^C v_E^C)$  and  $(3\sigma_E^S v_E^S, 3\sigma_E^C v_E^C)$ , for the environmental contribution, as detailed in Section 2.1. The variables  $\sigma_G^S, \sigma_G^C, \sigma_E^S$  and  $\sigma_E^C$  denote the standard deviation of the unnormalized shape and connectivity descriptors along the associated mode of co-variation. With a slight abuse of notation, in Figure 5-6, we denote these standard deviations by  $\sigma$ .

We then plot the shape and connectivity configurations identified by these tangent coordinates. Our main modes of genetic and environmental co-variation highlight an association between a global variation in brain shape and a global change in brain functional connectivity levels that seem to show a contrasting behavior in the genetic and environmental modes of co-variation.

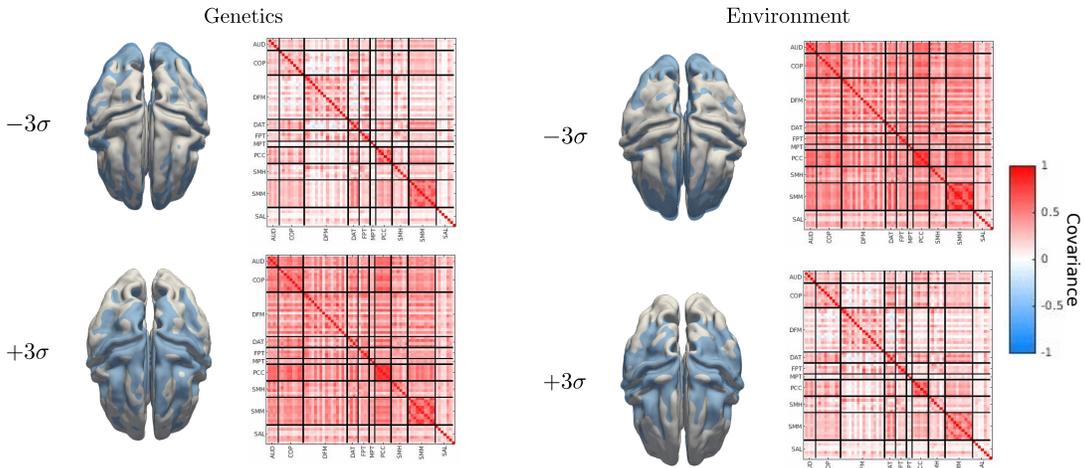


Fig 5: Illustration of the shape and connectivity CCA main modes of co-variation that are due to genetic and environmental contributions. Specifically, we display  $\pm 3\sigma$  changes in shape (gray-colored surfaces) and connectivity that are most correlated according to the estimated covariance structure  $\Sigma_G$  (left panel) and  $\Sigma_E$  (right panel). The variable  $\sigma$  here denotes the corresponding standard deviation of the unnormalized shape and connectivity descriptors in equation (4), once these have been projected along the directions representing the respective CCA modes of variation. Note that  $\sigma$  has different values for shape/connectivity and for genetic/environmental components. As in Figure 4, the blue surfaces are shown in order to have a fixed reference across the four panels and be able to appreciate the differences between the  $\pm 3\sigma$  gray surfaces. The main modes of co-variation display an association between a global change in shape and a global change in connectivity levels, with a contrasting behavior in the genetic and environmental components. Larger variations in connectivity levels are displayed between functional communities rather than within communities.

The proposed approach to modeling shape and connectivity allows us to capture both global and local variations. A more meticulous exploration of the results demonstrates a clear advantage of the adopted approach to shape modeling. In particular, in Figure 6, we show local shape changes that are associated with the mode of variation  $v_G^S$ . We can see non-trivial shape variations involving the formation of a sulcus. Capturing such fine-grained variations would in general not be possible with the simple shape descriptors (e.g., surface area) classically adopted in the neuroimaging literature.

In the supplementary material, we include more informative video plots that show the modes of co-variation as continuous shape and connectivity trajectories, associated with the tangent space directions  $(3t\sigma_G^S v_G^S, 3t\sigma_G^C v_G^C)$  and  $(3t\sigma_E^S v_E^S, 3t\sigma_E^C v_E^C)$ , for all  $t \in [-1, 1]$ . Moreover, we include the `.vtk` files of these trajectories, which allow for a more careful exploration of the results, by using a data visualization application, e.g., `Paraview` (Ayachit, 2015).

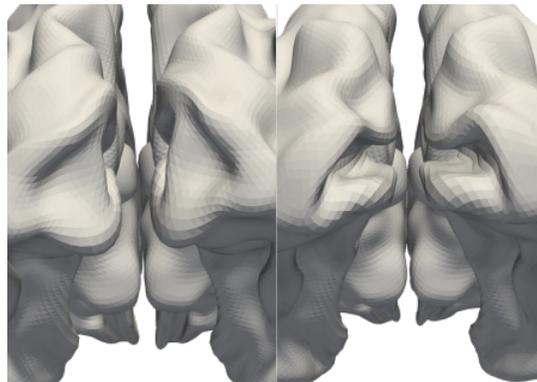


Fig 6: On the left, the shape configuration identified by the tangent vector  $-3\sigma v_G^S$ . On the right, the shape configuration identified by the tangent vector  $+3\sigma v_G^S$ . The tangent vector  $v_G^S$  represents shape changes, due to genetic contributions, in the main mode of co-variation between shape and connectivity. These are effectively a different view of the same gray surfaces shown on the top-left and bottom-left panels of Figure 5. The figure highlights the ability of the proposed framework to capture non-trivial localized variations in the brain shape, such as the formation of a sulcus.

**5. Discussion.** In this work, we propose a statistical Riemannian approach for the analysis of samples that are brain shapes and brain connectivity. In the proposed framework, we embed a variance component model that exploits family relatedness structure among samples to separate brain shape and connectivity co-variation that is due to genetic and environmental factors.

The proposed Riemannian modeling approach allows us to estimate trajectories in the spaces of shapes and connectivity that are constrained to belong to their respective non-Euclidean spaces of anatomically/physiologically meaningful estimates. Specifically, we are able to exclude shapes that are not topologically equivalent to the shapes in the sample, e.g., self-intersecting shapes, and we are able to exclude functional connectivity estimates that are not symmetric positive-definite objects. Moreover, the shape modeling approach proposed in this paper can also be easily extended to incorporate heterogeneous types of imaging data, such as volumetric representations of subcortical brain structures and bundles of axons estimated from Diffusion Tensor Imaging (Feydy et al., 2017). The proposed framework can be readily applied to the analysis of different anatomical objects.

When it comes to shapes and covariances, different representation models and metrics have been proposed in the literature. The choices we have made in this work are mainly driven by the reasons aforementioned. However, the exploration of different shape metrics, such as the Elastic Metric (Kurtek et al., 2011; Jermyn et al., 2012, 2017), is also a promising direction for future work. A particularly attractive property of the latter is the ability to integrate the registration step with the computation of the representation. Nonetheless, it is in principle more complicated to enforce topological constraints.

One limitation of the proposed variance component model is that it is based on a reduced dimension representation of shape and connectivity. It is, therefore, of interest to extend the current approach to work directly on the bivariate functional tangent space representations of shape and connectivity. Nevertheless, due to the high-dimensionality of the data, this is currently prohibitive, not only from a computational perspective but also due to the need to incorporate regularizing penalties to control for the nearly co-linear modes of co-variation that arise when estimating canonical correlation components from functional data. Further, in our analysis, we only focus on two factors contributing to the covariance structure of the traits: additive genetic and unique environmental factors. It is of course of interest to extend the proposed model to other factors, such as gene-environment interactions. It is also equally important to apply the proposed framework to the large imaging datasets nowadays available, such as the UK Biobank (Sudlow et al., 2015), while exploiting the genomic component of these datasets to estimate genetic relatedness. These datasets are more representative of the general population and could allow us to answer questions related to the effect of genes, age, and diseases on anatomical and functional properties of the brain.

## APPENDIX A: SIMULATIONS

In this section, we perform simulations to assess the finite sample estimation properties of the variance component model (1).

We generate two  $p \times p$  correlation matrices  $\Sigma_G$  and  $\Sigma_E$  as independent samples of a random correlation matrix which has a uniform distribution over the space of positive-definite correlation matrices (Joe, 2006). In order to obtain a simulation setting that is most similar to that of our application, we use the  $n \times n$  kinship matrix  $K_n$  from the HCP dataset, with  $n = 1001$ . Nevertheless, to study the effect of a hypothetical increase in sample size on our estimates, we introduce the  $nd \times nd$  kinship matrix  $I_d \otimes K_n$ . Such a kinship matrix describes a situation where  $nd$  samples are instead available, with  $d$  unrelated groups of  $n$  samples that have correlation structure represented by  $K_n$ .

We then generate our data following Model (1), i.e.,

$$\begin{aligned} A | G, E &= G + E \\ G &\sim \text{MVN}(0_{nd \times p}, I_d \otimes K_n, \Sigma_G), \\ E &\sim \text{MVN}(0_{nd \times p}, I_{nd}, \Sigma_E), \end{aligned}$$

where no fixed effects is included to reflect the setting of our final application. The  $nd \times p$  matrix  $A$  represents a set of simulated tangent space coordinates.

For every choice of  $d = 1, \dots, 5$ , and  $p = 3, 4$ , we generate 100 datasets and then apply the mixed effects model in Section 3 to compute the estimates  $\hat{\Sigma}_G$  and  $\hat{\Sigma}_E$  of  $\Sigma_G$  and  $\Sigma_E$ . We measure the accuracy of the estimate of the genetic component as  $\|\hat{\Sigma}_G - \Sigma_G\|_F$  and that of the environmental component as  $\|\hat{\Sigma}_E - \Sigma_E\|_F$ , where  $\|\cdot\|_F$  is the Frobenius norm.

We show the results of the simulations in Figure 7, where we can see that, as expected, an increase in sample size results in a lower estimation error. Moreover, as in the data generation model, we generate 1000 correlation matrices representing correlation matrices to be

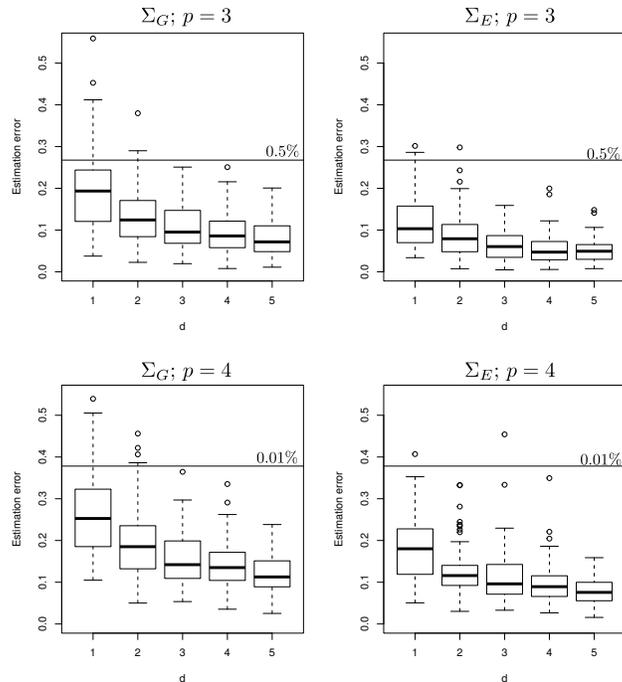


Fig 7: Boxplots describing the performances of the variance component model, in estimating  $\Sigma_G$  and  $\Sigma_E$ , as a function of the variable  $d = 1, \dots, 5$ , which is a multiplicative factor in the sample size  $nd$ . The estimation error is measured with the Frobenius norm of the difference between the covariance to be estimated and its estimate. We run 100 simulations for each choice of  $d$  and  $p$ . The horizontal lines denote the 0.5 and 0.01 percentiles, respectively for  $p = 3$  and  $p = 4$ , of the smallest estimation errors of a random guess estimator.

estimated. For each of the 1000 correlation matrices, we generate 1000 associated naive estimates, which are correlation matrices from the same distribution, and compute their Frobenius distance from the true correlation. The empirical distribution of the computed Frobenius norms describes the performance of the naive random estimator. We then select the 0.5 percentile, for  $p = 3$ , and 0.01 percentile, for  $p = 4$ , of the Frobenius norms. These are the horizontal lines in Figure 7. We can see that most of the estimates from the mixed-effects model are well below the selected threshold.

**Supplementary Material.** In the supplementary material, we provide animations that show the estimated modes of co-variation as continuous shape and connectivity trajectories. We also include the `.vtk` files of these trajectories, which can be visualized using, for instance, `Paraview`.

**Acknowledgments.** We wish to thank the editors and referees for the valuable comments and references.

**Funding.** JA wishes to gratefully acknowledge funding from Engineering and Physical Sciences Research Council (UK) EP/T017961/1.

## REFERENCES

ALMASY, L., DYER, T. D. and BLANGERO, J. (1997). Bivariate quantitative trait linkage analysis: Pleiotropy versus co-incident linkages. In *Genetic Epidemiology* **14** 953–958.

- AMOS, C. I. (1994). Robust variance-components approach for assessing genetic linkage in pedigrees. *American Journal of Human Genetics* **54** 535–543.
- ARGUILLÈRE, S., MILLER, M. I. and YOUNES, L. (2016). Diffeomorphic Surface Registration with Atrophy Constraints. *SIAM Journal on Imaging Sciences* **9** 975–1003.
- ARSIGNY, V., FILLARD, P., PENNEC, X. and AYACHE, N. (2006). Geometric means in a novel vector space structure on symmetric positive-definite matrices. *SIAM Journal on Matrix Analysis and Applications* **29** 328–347.
- ARSLAN, S., KTEA, S. I., MAKROPOULOS, A., ROBINSON, E. C., RUECKERT, D. and PARISOT, S. (2018). Human brain mapping: A systematic comparison of parcellation methods for the human cerebral cortex. *NeuroImage* **170** 5–30.
- AYACHIT, U. (2015). *The ParaView Guide: A Parallel Visualization Application*. Kitware, Inc., Clifton Park, NY, USA.
- BARBER, A. D., HEGARTY, C. E., LINDQUIST, M. and KARLSGODT, K. H. (2021). Heritability of Functional Connectivity in Resting State: Assessment of the Dynamic Mean, Dynamic Variance, and Static Connectivity across Networks. *Cerebral Cortex* **31** 2834–2844.
- BATES, D., MÄCHLER, M., BOLKER, B. M. and WALKER, S. C. (2015). Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* **67** 1–48.
- BULLMORE, E. and SPORNS, O. (2009). Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience* **10** 186–198.
- CALHOUN, V. D., LIU, J. and ADALI, T. (2009). A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data. *NeuroImage* **45** S163–S172.
- CARMICHAEL, I., CALHOUN, B. C., HOADLEY, K. A., TROESTER, M. A., GERADTS, J., COUTURE, H. D., OLSSON, L., PEROU, C. M., NIETHAMMER, M., HANNIG, J. and MARRON, J. S. (2019). Joint and individual analysis of breast cancer histologic images and genomic covariates.
- CHARLIER, B., CHARON, N. and TROUVÉ, A. (2017). The Fshape Framework for the Variability Analysis of Functional Shapes. *Foundations of Computational Mathematics* **17** 287–357.
- CHARLIER, B., NARDI, G. and TROUVÉ, A. (2015). The matching problem between functional shapes via a BV-penalty term: a  $\Gamma$ -convergence result. 1–31.
- CHARON, N. and TROUVÉ, A. (2014). Functional Currents: A New Mathematical Tool to Model and Analyse Functional Shapes. *Journal of Mathematical Imaging and Vision* **48** 413–431.
- CHEN, H. and WANG, Y. (2008). A Penalized spline approach to functional mixed effects model analysis. *Biometrics* **64** 751–761.
- DAHL, A., IOTCHKOVA, V., BAUD, A., JOHANSSON, Å., GYLLENSTEN, U., SORANZO, N., MOTT, R., KRANIS, A. and MARCHINI, J. (2016). A multiple-phenotype imputation method for genetic studies. *Nature Genetics* **48** 466–472.
- DESIKAN, R. S., SÉGONNE, F., FISCHL, B., QUINN, B. T., DICKERSON, B. C., BLACKER, D., BUCKNER, R. L., DALE, A. M., MAGUIRE, R. P., HYMAN, B. T., ALBERT, M. S. and KILLIANY, R. J. (2006). An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *NeuroImage* **31** 968–980.
- DRYDEN, I. L., KOLOYDENKO, A. and ZHOU, D. (2009). Non-Euclidean statistics for covariance matrices, with applications to diffusion tensor imaging. *The Annals of Applied Statistics* **3** 1102–1123.
- DRYDEN, I. L. and MARDIA, K. V. (2016). *Statistical Shape Analysis, with Applications in R. Wiley Series in Probability and Statistics*. John Wiley & Sons, Ltd, Chichester, UK.
- ELTZNER, B., HUCKEMANN, S. and MARDIA, K. V. (2018). Torus principal component analysis with applications to RNA structure. *The Annals of Applied Statistics* **12** 1332–1359.
- FENG, Q., JIANG, M., HANNIG, J. and MARRON, J. S. (2018). Angle-based joint and individual variation explained. *Journal of Multivariate Analysis* **166** 241–265.
- FEYDY, J., CHARLIER, B., VIALARD, F. X. and PEYRÉ, G. (2017). Optimal transport for diffeomorphic registration. In *Medical Image Computing and Computer Assisted Intervention - MICCAI 2017. MICCAI 2017. Lecture Notes in Computer Science* (M. DESCOTEAUX, L. MAIER-HEIN, A. FRANZ, P. JANNIN, D. COLLINS and S. DUCHESNE, eds.) **10433** 291–299. Springer, Cham.
- FIECAS, M., CRIBBEN, I., BAHKTIARI, R. and CUMMINE, J. (2017). A variance components model for statistical inference on functional connectivity networks. *NeuroImage* **149** 256–266.
- FISCHL, B., SERENO, M. I. and DALE, A. M. (1999). Cortical Surface-Based Analysis II: Inflation, Flattening, and a Surface-Based Coordinate System. *NeuroImage* **9** 195–207.
- FISCHL, B., VAN DER KOUWE, A., DESTRIEUX, C., HALGREN, E., SÉGONNE, F., SALAT, D. H., BUSA, E., SEIDMAN, L. J., GOLDSTEIN, J., KENNEDY, D., CAVINESS, V., MAKRIS, N., ROSEN, B. and DALE, A. M. (2004). Automatically Parcellating the Human Cerebral Cortex. *Cerebral Cortex* **14** 11–22.

- FRANCO, A. R., LING, J., CAPRIHAN, A., CALHOUN, V. D., JUNG, R. E., HEILEMAN, G. L. and MAYER, A. R. (2008). Multimodal and multi-tissue measures of connectivity revealed by joint independent component analysis. *IEEE Journal on Selected Topics in Signal Processing* **2** 986–997.
- GILMORE, J. H., KNICKMEYER, R. C. and GAO, W. (2018). Imaging structural and functional brain development in early childhood. *Nature Reviews Neuroscience* **19** 123–137.
- GLASSER, M. F., SOTIROPOULOS, S. N., WILSON, J. A., COALSON, T. S., FISCHL, B., ANDERSSON, J. L., XU, J., JBABDI, S., WEBSTER, M., POLIMENI, J. R., VAN ESSEN, D. C. and JENKINSON, M. (2013). The minimal preprocessing pipelines for the Human Connectome Project. *NeuroImage* **80** 105–124.
- GLASSER, M. F., SMITH, S. M., MARCUS, D. S., ANDERSSON, J. L. R., AUERBACH, E. J., BEHRENS, T. E. J., COALSON, T. S., HARMS, M. P., JENKINSON, M., MOELLER, S., ROBINSON, E. C., SOTIROPOULOS, S. N., XU, J., YACOB, E., UGURBIL, K. and VAN ESSEN, D. C. (2016). The Human Connectome Project’s neuroimaging approach. *Nature Neuroscience* **19** 1175–1187.
- GORDON, E. M., LAUMANN, T. O., ADEYEMO, B., HUCKINS, J. F., KELLEY, W. M. and PETERSEN, S. E. (2016). Generation and Evaluation of a Cortical Area Parcellation from Resting-State Correlations. *Cerebral Cortex* **26** 288–303.
- GUO, W. (2002). Functional Mixed Effects Models. *Biometrics* **58** 121–128.
- HAZLETT, H. C., GU, H., MUNSELL, B. C., KIM, S. H., STYNER, M., WOLFF, J. J., ELISON, J. T., SWANSON, M. R., ZHU, H., BOTTERON, K. N., COLLINS, D. L., CONSTANTINO, J. N., DAGER, S. R., ESTES, A. M., EVANS, A. C., FONOV, V. S., GERIG, G., KOSTOPOULOS, P., MCKINSTRY, R. C., PANDEY, J., PATERSON, S., PRUETT, J. R., SCHULTZ, R. T., SHAW, D. W., ZWAIGENBAUM, L. and PIVEN, J. (2017). Early brain development in infants at high risk for autism spectrum disorder. *Nature* **542** 348–351.
- IM, K., LEE, J. M., LYTTELTON, O., KIM, S. H., EVANS, A. C. and KIM, S. I. (2008). Brain size and cortical structure in the adult human brain. *Cerebral Cortex* **18** 2181–2191.
- JERMYN, I. H., KURTEK, S., KLASSEN, E. and SRIVASTAVA, A. (2012). Elastic shape matching of parameterized surfaces using square root normal fields. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)* **7576 LNCS** 804–817.
- JERMYN, I. H., KURTEK, S., LAGA, H. and SRIVASTAVA, A. (2017). Elastic Shape Analysis of Three-Dimensional Objects. *Synthesis Lectures on Computer Vision* **7** 1–185.
- JOE, H. (2006). Generating random correlation matrices based on partial correlations. *Journal of Multivariate Analysis* **97** 2177–2189.
- KANG, H. M., SUL, J. H., SERVICE, S. K., ZAITLEN, N. A., KONG, S. Y., FREIMER, N. B., SABATTI, C. and ESKIN, E. (2010). Variance component model to account for sample structure in genome-wide association studies. *Nature Genetics* **42** 348–354.
- KIRKPATRICK, M. and HECKMAN, N. (1989). A quantitative genetic model for growth, shape, reaction norms, and other infinite-dimensional characters. *Journal of Mathematical Biology* **27** 429–450.
- KURTEK, S., KLASSEN, E., DING, Z., JACOBSON, S. W., JACOBSON, J. L., AVISON, M. J. and SRIVASTAVA, A. (2011). Parameterization-invariant shape comparisons of anatomical surfaces. *IEEE Transactions on Medical Imaging* **30** 849–858.
- KURTEK, S., SRIVASTAVA, A., KLASSEN, E. and DING, Z. (2012). Statistical Modeling of Curves Using Shapes and Related Features. *Journal of the American Statistical Association* **107** 1152–1165.
- LANGE, K. (2002). *Mathematical and Statistical Methods for Genetic Analysis. Statistics for Biology and Health*. Springer New York, New York, NY.
- LENROOT, R. K. and GIEDD, J. N. (2008). The changing impact of genes and environment on brain development during childhood and adolescence: Initial findings from a neuroimaging study of pediatric twins. *Development and Psychopathology* **20** 1161–1175.
- LEWIS, A. S. and OVERTON, M. L. (2013). Nonsmooth optimization via quasi-Newton methods. *Mathematical Programming* **141** 135–163.
- LILA, E. and ASTON, J. A. D. (2020). Statistical Analysis of Functions on Surfaces, With an Application to Medical Imaging. *Journal of the American Statistical Association* **115** 1420–1434.
- LIU, B., WANG, L. and CAO, J. (2017). Estimating functional linear mixed-effects regression models. *Computational Statistics & Data Analysis* **106** 153–164.
- LUO, S., SONG, R., STYNER, M., GILMORE, J. H. and ZHU, H. (2019). FSEM: Functional Structural Equation Models for Twin Functional Data. *Journal of the American Statistical Association* **114** 344–357.
- MARDIA, K. V. and PATRANGENARU, V. (2005). Directions and projective shapes. *The Annals of Statistics* **33** 1666–1699.
- MARRON, J. S. and ALONSO, A. M. (2014). Overview of object oriented data analysis. *Biometrical Journal* **56** 732–753.
- MILLER, M. I., TROUVÉ, A. and YOUNES, L. (2006). Geodesic Shooting for Computational Anatomy. *Journal of Mathematical Imaging and Vision* **24** 209–228.

- MORRIS, J. S. and CARROLL, R. J. (2006). Wavelet-based functional mixed models. *Journal of the Royal Statistical Society. Series B: Statistical Methodology* **68** 179–199.
- PENNEC, X., FILLARD, P. and AYACHE, N. (2006). A riemannian framework for tensor computing. *International Journal of Computer Vision* **66** 41–66.
- PIGOLI, D., ASTON, J. A. D., DRYDEN, I. L. and SECCHI, P. (2014). Distances and inference for covariance operators. *Biometrika* **101** 409–422.
- PIZER, S. M., JUNG, S., GOSWAMI, D., VICORY, J., ZHAO, X., CHAUDHURI, R., DAMON, J. N., HUCKEMANN, S. and MARRON, J. S. (2013). Nested Sphere Statistics of Skeletal Models. In *Innovations for Shape Analysis. Mathematics and Visualization*. (M. Breuß, A. Bruckstein and P. Maragos, eds.) 93–115. Springer, Berlin, Heidelberg.
- PIZZAGALLI, F., AUZIAS, G., YANG, Q., MATHIAS, S. R., FASKOWITZ, J., BOYD, J. D., AMINI, A., RIVIÈRE, D., MCMAHON, K. L., DE ZUBICARAY, G. I., MARTIN, N. G., MANGIN, J. F., GLAHN, D. C., BLANGERO, J., WRIGHT, M. J., THOMPSON, P. M., KOCHUNOV, P. and JAHANSHAD, N. (2020). The reliability and heritability of cortical folds and their genetic correlations across hemispheres. *Communications Biology* **3**.
- POWER, J. D., COHEN, A. L., NELSON, S. M., WIG, G. S., BARNES, K. A., CHURCH, J. A., VOGEL, A. C., LAUMANN, T. O., MIEZIN, F. M., SCHLAGGAR, B. L. and PETERSEN, S. E. (2011). Functional Network Organization of the Human Brain. *Neuron* **72** 665–678.
- QIN, L. (2005). Functional mixed-effects model for periodic data. *Biostatistics* **7** 225–234.
- RAMSAY, J. and SILVERMAN, W. B. (2005). *Functional Data Analysis. Springer Series in Statistics*. Springer-Verlag, New York.
- REIMHERR, M. and NICOLAE, D. (2016). Estimating Variance Components in Functional Linear Models With Applications to Genetic Heritability. *Journal of the American Statistical Association* **111** 407–422.
- RISK, B. B. and ZHU, H. (2019). ACE of space: estimating genetic components of high-dimensional imaging data. *Biostatistics* 1–45.
- ROBINSON, E. C., JBABDI, S., GLASSER, M. F., ANDERSSON, J., BURGESS, G. C., HARMS, M. P., SMITH, S. M., VAN ESSEN, D. C. and JENKINSON, M. (2014). MSM: A new flexible framework for Multimodal Surface Matching. *NeuroImage* **100** 414–426.
- ROBINSON, E. C., GARCIA, K., GLASSER, M. F., CHEN, Z., COALSON, T. S., MAKROPOULOS, A., BOZEK, J., WRIGHT, R., SCHUH, A., WEBSTER, M., HUTTER, J., PRICE, A., CORDERO GRANDE, L., HUGHES, E., TUSOR, N., BAYLY, P. V., VAN ESSEN, D. C., SMITH, S. M., EDWARDS, A. D., HAJNAL, J., JENKINSON, M., GLOCKER, B. and RUECKERT, D. (2018). Multimodal surface matching with higher-order smoothness constraints. *NeuroImage* **167** 453–465.
- SCHEIPL, F., STAIUCU, A. M. and GREVEN, S. (2015). Functional Additive Mixed Models. *Journal of Computational and Graphical Statistics* **24** 477–501.
- SHI, M., WEISS, R. E. and TAYLOR, J. M. G. (1996). An Analysis of Paediatric CD4 Counts for Acquired Immune Deficiency Syndrome Using Flexible Random Curves. *Journal of the Royal Statistical Society. Series C (Applied Statistics)* **45** 151.
- SMITH, S. M., BECKMANN, C. F., ANDERSSON, J., AUERBACH, E. J., BIJSTERBOSCH, J., DOUAUD, G., DUFF, E., FEINBERG, D. A., GRIFFANTI, L., HARMS, M. P., KELLY, M., LAUMANN, T., MILLER, K. L., MOELLER, S., PETERSEN, S., POWER, J., SALIMI-KHORSHIDI, G., SNYDER, A. Z., VU, A. T., WOOLRICH, M. W., XU, J., YACOUB, E., UĞURBIL, K., VAN ESSEN, D. C. and GLASSER, M. F. (2013). Resting-state fMRI in the Human Connectome Project. *NeuroImage* **80** 144–168.
- SMITH, S. M., NICHOLS, T. E., VIDAURRE, D., WINKLER, A. M., BEHRENS, T. E. J., GLASSER, M. F., UĞURBIL, K., BARCH, D. M., VAN ESSEN, D. C. and MILLER, K. L. (2015). A positive-negative mode of population covariation links brain connectivity, demographics and behavior. *Nature Neuroscience* **18** 1565–1567.
- SU, J., KURTEK, S., KLASSEN, E. and SRIVASTAVA, A. (2014). Statistical analysis of trajectories on Riemannian manifolds: Bird migration, hurricane tracking and video surveillance. *The Annals of Applied Statistics* **8** 530–552.
- SUDLOW, C., GALLACHER, J., ALLEN, N., BERAL, V., BURTON, P., DANESH, J., DOWNEY, P., ELLIOTT, P., GREEN, J., LANDRAY, M., LIU, B., MATTHEWS, P., ONG, G., PELL, J., SILMAN, A., YOUNG, A., SPROSEN, T., PEAKMAN, T. and COLLINS, R. (2015). UK Biobank: An Open Access Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age. *PLoS Medicine* **12** 1001779.
- SUI, J., PEARLSON, G., CAPRIHAN, A., ADALI, T., KIEHL, K. A., LIU, J., YAMAMOTO, J. and CALHOUN, V. D. (2011). Discriminating schizophrenia and bipolar disorder by fusing fMRI and DTI in a multimodal CCA+ joint ICA model. *NeuroImage* **57** 839–855.
- VAILLANT, M., MILLER, M. I., YOUNES, L. and TROUVÉ, A. (2004). Statistics on diffeomorphisms via tangent space representations. *NeuroImage* **23** S161–S169.

- VAN ESSEN, D. C., GLASSER, M. F., DIERKER, D. L., HARWELL, J. and COALSON, T. (2012). Parcellations and Hemispheric Asymmetries of Human Cerebral Cortex Analyzed on Surface-Based Atlases. *Cerebral Cortex* **22** 2241–2262.
- VAN ESSEN, D. C., SMITH, S. M., BARCH, D. M., BEHRENS, T. E. J., YACOUB, E. and UGURBIL, K. (2013). The WU-Minn Human Connectome Project: An overview. *NeuroImage* **80** 62–79.
- WANG, B., SVERDLOV, S. and THOMPSON, E. (2017). Efficient Estimation of Realized Kinship from Single Nucleotide Polymorphism Genotypes. *Genetics* **205** 1063–1078.
- WIG, G. S., LAUMANN, T. O. and PETERSEN, S. E. (2014). An approach for parcellating human cortical areas using resting-state correlations. *NeuroImage* **93** 276–291.
- WU, H. and ZHANG, J.-T. (2002). Local Polynomial Mixed-Effects Models for Longitudinal Data. *Journal of the American Statistical Association* **97** 883–897.
- XIA, C. H., MA, Z., CIRIC, R., GU, S., BETZEL, R. F., KACZKURKIN, A. N., CALKINS, M. E., COOK, P. A., GARCÍA DE LA GARZA, A., VANDEKAR, S. N., CUI, Z., MOORE, T. M., ROALF, D. R., RUPAREL, K., WOLF, D. H., DAVATZIKOS, C., GUR, R. C., GUR, R. E., SHINOHARA, R. T., BASSETT, D. S. and SATTERTHWAITE, T. D. (2018). Linked dimensions of psychopathology and connectivity in functional brain networks. *Nature Communications* **9**.
- XUE, W., DUBOIS BOWMAN, F., PILEGGI, A. V. and MAYER, A. R. (2015). A multimodal approach for determining brain networks by jointly modeling functional and structural connectivity. *Frontiers in Computational Neuroscience* **9** 1–11.
- YEO, B. T. T., SABUNCU, M. R., VERCAUTEREN, T., AYACHE, N., FISCHL, B. and GOLLAND, P. (2010). Spherical Demons: Fast Diffeomorphic Landmark-Free Surface Registration. *IEEE Transactions on Medical Imaging* **29** 650–668.
- YEO, B. T. T., KRIENEN, F. M., SEPULCRE, J., SABUNCU, M. R., LASHKARI, D., HOLLINSHEAD, M., ROFFMAN, J. L., SMOLLER, J. W., ZÖLLEI, L., POLIMENI, J. R., FISCHL, B., LIU, H. and BUCKNER, R. L. (2011). The organization of the human cerebral cortex estimated by intrinsic functional connectivity. *Journal of Neurophysiology* **106** 1125–1165.
- YOUNES, L. (2010). *Shapes and Diffeomorphisms. Applied Mathematical Sciences*. Springer, Berlin, Heidelberg.
- ZHANG, Z., ALLEN, G. I., ZHU, H. and DUNSON, D. (2019). Tensor network factorizations: Relationships between brain structural connectomes and traits. *NeuroImage* **197** 330–343.
- ZHOU, L., HUANG, J. Z. and CARROLL, R. J. (2008). Joint modelling of paired sparse functional data using principal components. *Biometrika* **95** 601–619.
- ZHOU, X. and STEPHENS, M. (2014). Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nature Methods* **11** 407–409.
- ZITOVÁ, B. and FLUSSER, J. (2003). Image registration methods: a survey. *Image and Vision Computing* **21** 977–1000.