

1 Brain kernel: a new spatial covariance function for fMRI data

2
3 Anqi Wu^{1†}, Samuel A. Nastase^{2,3}, Christopher A. Baldassano⁴, Nicholas B. Turk-Browne⁵, Kenneth A.
4 Norman^{2,3}, Barbara E. Engelhardt⁶, Jonathan W. Pillow^{2,3}

5 **1** Center for Theoretical Neuroscience, Columbia University, New York City, NY, USA

6 **2** Princeton Neuroscience Institute, Princeton University, Princeton, NJ, USA

7 **3** Department of Psychology, Princeton University, Princeton, NJ, USA

8 **4** Department of Psychology, Columbia University, New York City, NY, USA

9 **5** Department of Psychology, Yale University, New Haven, CT, USA

10 **6** Department of Computer Science, Princeton University, Princeton, NJ, USA

11 † Correspondence email: anqiwu.angela@gmail.com

12 Abstract

13 A key problem in functional magnetic resonance imaging (fMRI) is to estimate spatial activity patterns
14 from noisy high-dimensional signals. Spatial smoothing provides one approach to regularizing such
15 estimates. However, standard smoothing methods ignore the fact that correlations in neural activity may
16 fall off at different rates in different brain areas, or exhibit discontinuities across anatomical or functional
17 boundaries. Moreover, such methods do not exploit the fact that widely separated brain regions may exhibit
18 strong correlations due to bilateral symmetry or the network organization of brain regions. To capture
19 this non-stationary spatial correlation structure, we introduce the *brain kernel*, a continuous covariance
20 function for whole-brain activity patterns. We define the brain kernel in terms of a continuous nonlinear
21 mapping from 3D brain coordinates to a latent embedding space, parametrized with a Gaussian process
22 (GP). The brain kernel specifies the prior covariance between voxels as a function of the distance between
23 their locations in embedding space. The GP mapping warps the brain nonlinearly so that highly correlated
24 voxels are close together in latent space, and uncorrelated voxels are far apart. We estimate the brain
25 kernel using resting-state fMRI data, and we develop an exact, scalable inference method based on
26 block coordinate descent to overcome the challenges of high dimensionality (10-100K voxels). Finally, we
27 illustrate the brain kernel's usefulness with applications to brain decoding and factor analysis with multiple
28 task-based fMRI datasets.

29 **Keywords:** Brain kernel, Gaussian process, latent variable model, brain decoding, factor modeling, resting-
30 state fMRI, task fMRI.

31 1 Introduction

32 An important problem in neuroscience is to characterize the covariance of high-dimensional neural activity.
33 Understanding this covariance structure could provide insight into the brain's functional organization and
34 help regularize estimates of encoding or decoding models. Although advances have been made in both
35 theory and methodology for estimating large covariance [1, 2] and precision matrices [3, 4], few methods
36 have been designed with the particular challenges of functional magnetic resonance imaging (fMRI) data
37 in mind (see [5] for an exception).

38 One of the challenges of modeling the covariance of fMRI data is that the spatial discretization of the
39 brain may differ across experiments. fMRI measures blood oxygenation level dependent (BOLD) signals
40 in discrete spatial regions called "voxels". Each voxel represents a tiny cube of brain tissue. Although
41 brains are typically registered to an anatomical template in a standard space, such as the volumetric

42 Montreal Neurological Institute (MNI) template or the surface-based template used by HCP [6], these
43 spaces differ in resolution and geometry [7]. In many cases, brains are aligned onto the “same” 3D space
44 but with different voxel coordinates. A covariance matrix for one set of voxels cannot be applied to data
45 registered to a different set of voxels. Thus, modeling the covariance of fMRI data presently requires a
46 new covariance matrix to be constructed whenever a different set of voxels is used.

47 A second challenge for fMRI covariance estimation is spatial nonstationarity. Standard spatial smoothing
48 models assume that correlation falls off as a function of the Euclidean distance between voxels. In real
49 brains, however, correlation patterns depend on relationships to anatomical and functional boundaries,
50 and may exhibit strong dependencies over long distances due to bilateral symmetry and the network
51 organization of brain regions.

52 To address these challenges, we propose the *brain kernel*, a continuous covariance function for whole-
53 brain fMRI data. This function arises from a generative model of fMRI data, and seeks to describe
54 covariance of neural signals across the entire brain [8]. Specifically, the brain kernel defines, for any
55 finite set of n voxel locations in the brain, a positive definite $n \times n$ covariance matrix over n -dimensional
56 vectors of neural activity at those locations. The brain kernel improves upon prior work by i) capturing fMRI
57 voxel covariance matrices for any registration reference, and ii) capturing spatial nonstationarity through a
58 nonlinear latent manifold.

59 Our approach uses a Gaussian process to parametrize a continuous nonlinear mapping from 3D brain
60 coordinates to a latent embedding space, such that correlations in neural activity fall off as a fixed function
61 of distance in the latent space. Thus, the nonlinear function seeks to warp the 3D brain in order to place
62 locations with correlated neural activity at nearby locations in the latent space. Locations with uncorrelated
63 activity, conversely, are mapped to more distant points in the latent space, even if they are physically close
64 together in the brain.

65 The paper is organized as follows. In Sec. 2 we provide a brief overview of Gaussian process models. In
66 Sec. 3, we formally introduce the brain kernel model for fMRI data. In Sec. 4, we describe an efficient
67 inference method for fitting the brain kernel, and illustrate the challenges and benefits of an exact inference
68 method using simulated and real fMRI datasets. In Sec. 5, we describe the brain kernel fit to whole-brain
69 resting-state fMRI data. Finally, in Sec. 6, we demonstrate the usefulness of the inferred brain kernel with
70 applications to decoding and factor modeling.

71 **2 Mathematical background**

72 Before introducing the brain kernel model, we briefly review the mathematical building blocks for Gaussian
73 process models.

74 **2.1 Gaussian processes (GPs)**

75 Gaussian processes provide a flexible and tractable prior distribution over nonlinear functions [9]. A GP
76 is parametrized by a mean function $m(\mathbf{x})$, which specifies the mean value of the function $f(\mathbf{x})$ at input
77 point \mathbf{x} , and a covariance function $k(\mathbf{x}, \mathbf{x}')$, which specifies $\text{cov}(f(\mathbf{x}), f(\mathbf{x}'))$, the covariance between the
78 function values $f(\mathbf{x})$ and $f(\mathbf{x}')$, for any pair of inputs \mathbf{x} and \mathbf{x}' .

79 Technically a GP is a random process for which the values taken at any finite set of input points has a
80 well-defined multivariate Gaussian distribution. The mean and covariance of that Gaussian are given by
81 evaluating the mean and covariance functions at the corresponding set of input points. Let $(\mathbf{x}_1, \dots, \mathbf{x}_n)$
82 denote a collection of n points in the input domain. If a function f has a Gaussian process distribution,
83 $f \sim \mathcal{GP}(m, k)$, then the vector of function values $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^\top$ has a multivariate Gaussian

84 distribution:

$$\mathbf{f} \sim \mathcal{N}(\mathbf{m}, \mathbf{K}), \quad (1)$$

85 where $\mathbf{m} = (m(\mathbf{x}_1), \dots, m(\mathbf{x}_n))^T$ is the mean vector, and \mathbf{K} is the $(n \times n)$ covariance matrix whose
86 i, j 'th element is $k(\mathbf{x}_i, \mathbf{x}_j)$.

87 2.2 GP regression

88 A common application of GPs is to predict function values at test points given a set of training data
89 consisting of observed inputs and function values. In GP regression, these predictions come from the
90 conditional distribution over unknown function values given the observed values.

91 Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ denote a set of n input points and let $\mathbf{f} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))^T$ denote the
92 vector of function values observed at these points. Here we assume the observed data is noiseless,
93 and we will consider noisy observations in Sec. 3. We consider a set of n^* novel input points $\mathbf{X}_* =$
94 $(\mathbf{x}_{n+1}, \dots, \mathbf{x}_{n+n^*})$, for which we would like to predict the corresponding (unobserved) function values,
95 denoted $\mathbf{f}_* = (f(\mathbf{x}_{n+1}), \dots, f(\mathbf{x}_{n+n^*}))^T$.

96 The GP gives us the following joint prior distribution over the observed and unobserved function values:

$$\begin{bmatrix} \mathbf{f} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m} \\ \mathbf{m}_* \end{bmatrix}, \begin{bmatrix} \mathbf{K}_{nn} & \mathbf{K}_{n*} \\ \mathbf{K}_{*n} & \mathbf{K}_{**} \end{bmatrix} \right), \quad (2)$$

97 where $\mathbf{m}_* = (m(\mathbf{x}_{n+1}), \dots, m(\mathbf{x}_{n+n^*}))^T$ is the mean for novel test points in \mathbf{X}_* , and matrices \mathbf{K}_{n*} , \mathbf{K}_{*n} ,
98 and \mathbf{K}_{**} are of size $(n \times n^*)$, $(n^* \times n)$, and $(n^* \times n^*)$, respectively, formed by evaluating the covariance
99 function k at the relevant points in \mathbf{X} and \mathbf{X}_* .

By applying the standard formula for Gaussian conditional distributions, we obtain the following conditional
distribution of \mathbf{f}_* given \mathbf{f} :

$$\mathbf{f}_* | \mathbf{f} \sim \mathcal{N}(\mu(\mathbf{X}_*), \sigma^2(\mathbf{X}_*)), \quad (3)$$

where mean and covariance are given by

$$\mu(\mathbf{X}_*) = \mathbf{K}_{*n} \mathbf{K}_{nn}^{-1} (\mathbf{f} - \mathbf{m}), \quad (4)$$

$$\sigma^2(\mathbf{X}_*) = \mathbf{K}_{**} - \mathbf{K}_{*n} \mathbf{K}_{nn}^{-1} \mathbf{K}_{n*}. \quad (5)$$

100 GP regression uses eq. 4, the posterior mean of the function given the training data, to predict function
101 values at test points \mathbf{X}_* given the training data $\{\mathbf{X}, \mathbf{f}\}$.

102 Although we have assumed so far that the function f is scalar-valued, we can extend the GP regression
103 framework to vector-valued functions by using a separate GP for each output dimension of f .

104 2.3 Gaussian process latent variable model (GPLVM)

105 The Gaussian process latent variable model (GPLVM) [10] extends the GP regression model by treating the
106 inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ as latent variables, distributed independently with standard normal distributions:
107 $\mathbf{x}_i \sim \mathcal{N}(0, I)$, for $i = 1, \dots, n$. Unlike the regression setting, the goal of the GPLVM is to account for
108 structure of an observed dataset in terms of a set of locations in a latent space, which are transformed
109 by a smooth function to generate the observed function values. The primary quantity of interest in the
110 GPLVM is therefore $P(\mathbf{X} | \mathbf{f})$, the conditional distribution over the latent variables given the observations.

111 The smoothness induced by the GP prior distribution ensures a smooth mapping from latent variables to
 112 observed data. This means that when observed values f_i and f_j are close together, they are likely to arise
 113 from points x_i and x_j that are close together in the latent space. Conversely, if observed values f_i
 114 and f_j are far apart, smoothness of f implies that the corresponding input values lie far apart in the latent
 115 space.
 116 In the following, we will derive the brain kernel model, which extends the GPLVM to allow latent locations
 117 arise from a smooth nonlinear transformation of the 3D voxel locations in the brain. The brain kernel
 118 model can also be seen as a “deep Gaussian process” model [11] with two layers (see Appendix A for
 119 complete details).

120 3 The brain kernel model

121 Here we introduce the brain kernel (BK) model, which is a probabilistic model of fMRI measurements at
 122 an arbitrary set of 3D spatial voxel locations. The brain kernel itself is a covariance function for neural
 123 activity that arises under this BK model, which we will infer from registered fMRI data.

124 3.1 Nonlinear embedding function

125 The first component of the brain kernel model is a nonlinear function, $f : \mathbb{R}^3 \rightarrow \mathbb{R}^d$, which provides a
 126 continuous nonlinear mapping from 3D brain coordinates to a d -dimensional latent embedding space. The
 127 goal of this mapping is to embed brain regions with similar activity at nearby locations in the embedded
 128 space. Typically, we consider $d > 3$, so that the embedding is higher-dimensional than the three
 129 physical dimensions of the brain. This gives the embedding flexibility to capture complex non-smooth
 130 dependencies between brain regions. One example of such a dependency is the functional symmetry of
 131 the two hemispheres [12–14], which suggests that one might wish to map symmetric points on the two
 132 hemispheres to nearby points in the embedded latent space. This would not be possible with a continuous
 133 mapping in three dimensions. But, in four dimensions, one can fold the three-dimensional brain along the
 134 fourth dimension, analogous to the way that folding a 2D brain slice along the mid-line would allow for
 135 close alignment of symmetric points from the two hemispheres.

136 Let $\mathbf{x} \in \mathbb{R}^3$ denote an input vector, specifying the three-dimensional location of a voxel in the brain, and let
 137 $\mathbf{z} \in \mathbb{R}^d$ denote the output of f , so that $\mathbf{z} = f(\mathbf{x})$ is the d -dimensional embedding location of a voxel at
 138 \mathbf{x} . Thus, for a set of voxel locations $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$, the embedded locations in the latent space are
 139 $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n) = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$.

140 To impose smoothness on the embedding function, we place a GP prior on f . We use a linear mean
 141 function, $m(\mathbf{x}) = \mathbf{B}\mathbf{x}$, where \mathbf{B} is a $d \times 3$ matrix. This choice ensures that the embedding defaults to a
 142 linearly stretched version of the brain in the absence of likelihood terms. Because f is a vector function
 143 with outputs of dimension d , the mean function output is also d -dimensional. For the covariance function,
 144 we use a Gaussian or “radial basis function” (RBF) covariance for each output dimension of f :

$$k_f(\mathbf{x}_i, \mathbf{x}_j) = r \exp\left(-\frac{1}{2\delta^2} \|\mathbf{x}_i - \mathbf{x}_j\|_2^2\right). \quad (6)$$

145 The hyperparameters governing this covariance function consist of a marginal variance r and a length
 146 scale δ , which control the range and smoothness of f , respectively. The GP prior over each output
 147 dimension of f can therefore be written as

$$f_j(\cdot) \sim GP(\mathbf{b}_j, k_f), \quad (7)$$

148 where $f_j(\cdot)$ denotes the j th output dimension of the function $f(\cdot)$, and \mathbf{b}_j is the j th row of the matrix \mathbf{B} .
 149 The prior is therefore governed by a set of hyperparameters denoted $\theta_f = \{\mathbf{B}, r, \delta\}$. Thus, all output

150 dimensions of the function f are assumed *a priori* independent with the same covariance function and
151 differing mean functions.

152 This GP prior over the function f implies a multivariate normal prior over any set of embedded voxel
153 locations \mathbf{Z} . Let \mathbf{z}_j denote the j th latent embedding of the entire set of brain voxels in the training data.
154 Then the prior over \mathbf{z}_j given the true voxel locations \mathbf{X} is

$$p(\mathbf{z}_j | \mathbf{X}) = \mathcal{N}(\mathbf{b}_j \mathbf{X}, \mathbf{K}), \quad (8)$$

155 where $\mathbf{K} \in \mathbb{R}^{n \times n}$ is the covariance matrix with the i, j th element given by $k(\mathbf{x}_i, \mathbf{x}_j)$ (eq. 6).

156 3.2 From embedding space to neural activity

157 The second component of the brain kernel model is a probability distribution over neural activity as a
158 function of locations in embedding space. Our modeling assumption is that neural activity changes
159 smoothly as a function of locations in embedding space, or equivalently, that correlations in neural activity
160 decrease smoothly with distance in latent space. We formalize this assumption using the brain kernel,
161 which provides a mapping from latent embedding locations to a covariance matrix for neural activity.

162 Let $\mathbf{v} \in \mathbb{R}^n$ denote a vector of neural activity from n voxels with positions $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ and latent
163 embedding locations $\mathbf{Z} = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n))$. The BK model assumes this neural activity vector has a
164 multivariate Gaussian distribution with zero mean and covariance determined by the brain kernel:

$$\mathbf{v} \sim \mathcal{N}(0, \mathbf{C}_{nn}), \quad (9)$$

165 where

$$\mathbf{C}_{nn} = \begin{bmatrix} \kappa_{BK}(\mathbf{x}_1, \mathbf{x}_1) & \cdots & \kappa_{BK}(\mathbf{x}_1, \mathbf{x}_n) \\ \vdots & \ddots & \vdots \\ \kappa_{BK}(\mathbf{x}_n, \mathbf{x}_1) & \cdots & \kappa_{BK}(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix} \quad (10)$$

166 is the covariance matrix, which results from applying the brain kernel $\kappa_{BK}(\cdot, \cdot)$ to every pair of voxel
167 locations $(\mathbf{x}_i, \mathbf{x}_j)$ in the set \mathbf{X} .

The brain kernel itself is the bivariate function $\kappa_{BK} : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}$ from pairs of 3D voxel locations to a
covariance of neural activity at those pairs of locations:

$$\kappa_{BK}(\mathbf{x}_i, \mathbf{x}_j) = \rho \exp\left(-\frac{1}{2}\|\mathbf{f}(\mathbf{x}_i) - \mathbf{f}(\mathbf{x}_j)\|_2^2\right) = \rho \exp\left(-\frac{1}{2}\|\mathbf{z}_i - \mathbf{z}_j\|_2^2\right), \quad (11)$$

168 where ρ is the marginal variance. The brain kernel therefore specifies a positive semidefinite covariance
169 matrix for neural activity at any set of 3D voxel locations, which is a function of the embedded latent
170 locations of those voxels via the nonlinear function f . The brain kernel model transforms the data
171 representation from voxel space to the latent embedding space, and this transformation explicitly estimates
172 the covariance over neural activity (Figure 1).

173 3.3 From neural activity to BOLD signal

174 Next, we assume that the experimenter does not directly measure the neural activity vector \mathbf{v} , but instead
175 receives measurements corrupted by independent Gaussian noise. If \mathbf{y} denotes the vector of fMRI
176 measurements, we assume $y_i = v_i + \xi_i$, where i is the index for voxels and $\xi_i \sim \mathcal{N}(0, \sigma^2)$ represents
177 measurement noise. This induces the following marginal distribution over fMRI measurements given the
178 embedding:

$$\mathbf{y} | f \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{nn} + \sigma^2 I_n). \quad (12)$$

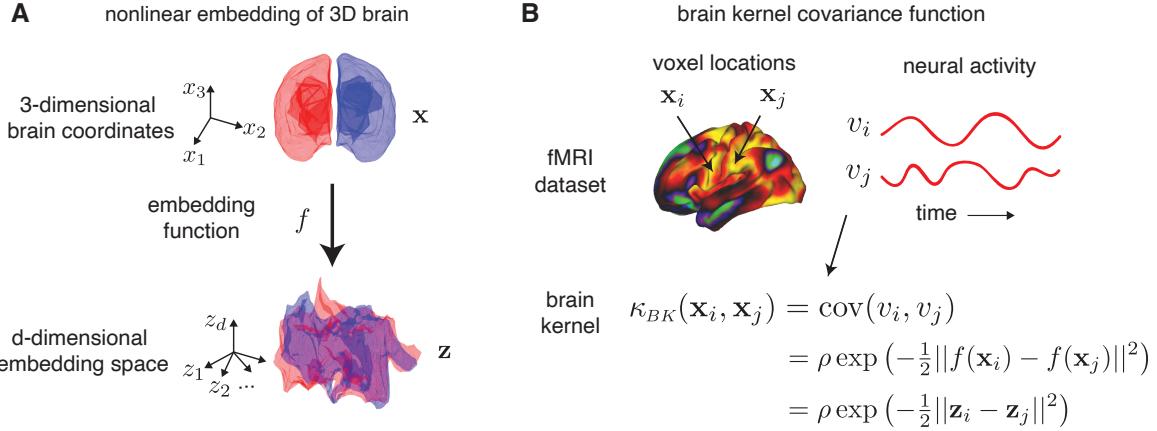


Fig 1. The diagram of the brain kernel model. **A.** The nonlinear latent embedding of the 3D coordinates into a d -dimensional latent space using a function f . This f function is sampled from a GP prior. **B.** Given the voxel locations and the BOLD activity of these two locations in an fMRI dataset (top right), our goal is to construct the brain kernel with the locations as inputs that matches their covariance of the BOLD activity (bottom right). The covariance is equal to the Euclidean-based kernel of the voxels' embedding in the latent space.

179 Note however that the brain kernel generates \mathbf{C}_{nn} , the covariance of the underlying neural activity \mathbf{v} , as
 180 opposed to the covariance of the noisy fMRI measurements \mathbf{y} ; the covariance of these measurements is
 181 $(\mathbf{C}_{nn} + \sigma^2 \mathbf{I}_n)$. For simplicity, we will omit the subscript in \mathbf{C}_{nn} in the following text and write simply \mathbf{C} .

182 The full set of hyperparameters governing the brain kernel model are therefore $\theta = \{\mathbf{B}, r, \delta, \rho, \sigma^2\}$, where
 183 $\{\mathbf{B}, r, \delta\}$ describe the nonlinear embedding function, ρ is the marginal variance of neural activity, and σ^2
 184 is the variance of additive Gaussian noise.

185 4 Inference methods

186 To fit the brain kernel model, we estimate the latent embeddings of all voxels \mathbf{Z} as well as the model
 187 hyperparameters $\theta = \{\mathbf{B}, r, \delta, \rho, \sigma^2\}$ given a time series of whole-brain fMRI measurements \mathbf{Y} as
 188 well as voxels' 3D locations \mathbf{X} . More specifically, $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_T) \in \mathbb{R}^{n \times T}$, where $\mathbf{y}_t \in \mathbb{R}^n$ refers
 189 to the vector of fMRI measurements at time index $t \in \{1, \dots, T\}$. Let $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n) \in \mathbb{R}^{3 \times n}$,
 190 where $\mathbf{x}_i \in \mathbb{R}^3$ and $i \in \{1, \dots, n\}$, denote the set of n 3D voxel locations for this dataset. Let $\mathbf{Z} =$
 191 $(\mathbf{z}_1, \dots, \mathbf{z}_n) = (f(\mathbf{x}_1), \dots, f(\mathbf{x}_n)) \in \mathbb{R}^{d \times n}$, where $\mathbf{z}_i \in \mathbb{R}^d$ denotes the set of n voxels' d -dimensional
 192 latent representations.

193 We propose two empirical estimators: penalized least squares (PLS) and maximum a posteriori (MAP).
 194 Briefly, PLS formulates an objective function by minimizing the squared error between the sample
 195 covariance of the data $\text{cov}(\mathbf{Y})$, and the model-defined covariance \mathbf{C} (eq. 10). MAP solves the problem
 196 using conventional Bayesian probabilistic inference. Since we have already defined the distribution to
 197 generate fMRI measurements from the latent embeddings (eq. 12) and the prior for the latent embeddings
 198 (eq. 8), we can estimate the latent embeddings \mathbf{Z} using MAP methods. The goal of both inference methods
 199 is to find the latent embedding locations \mathbf{Z} and the model hyperparameters such that the neural activity
 200 covariance \mathbf{C} resembles the sample covariance of \mathbf{Y} as closely as possible.

201 4.1 Penalized least squares (PLS) estimation

202 Estimating large covariance matrices is a fundamental problem in modern multivariate analysis. Re-
 203 searchers attempt to find an estimator that resembles the sample covariance while also satisfying structural
 204 assumptions about the data [15, 16]. In prior work, Fan et al. [16] showed that a generalized thresholding
 205 covariance estimator can be cast as a penalized least squares (PLS) problem:

$$\hat{\mathbf{C}} = \underset{\mathbf{C}}{\operatorname{argmin}} \{ \|\mathbf{S} - \mathbf{C}\| + R(\mathbf{C}) \}, \quad (13)$$

206 where $R(\cdot)$ is a penalty function that imposes structure on the covariance matrix \mathbf{C} , and \mathbf{S} is the sample
 207 covariance of measured neural activity, defined as

$$\mathbf{S} = \frac{1}{T-1} \sum_{t=1}^T (\mathbf{y}_t - \bar{\mathbf{y}})(\mathbf{y}_t - \bar{\mathbf{y}})^\top, \quad \bar{\mathbf{y}} = \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t. \quad (14)$$

Sparsity in \mathbf{C} is often encoded using a shrinkage penalty. However, we abandon sparsity in the brain kernel because the covariance of brain activity may have dense structures. Instead, we regularize the latent subspace of the covariance matrix using the normal prior on the latent embedding \mathbf{Z} (eq. 8), which is a Bayesian regularization of the log likelihood with the form $\operatorname{tr}[(\mathbf{Z} - \mathbf{BX})\mathbf{K}^{-1}(\mathbf{Z} - \mathbf{BX})^\top]$. This is equivalent to an l_2 -norm penalty on $\mathbf{K}^{-\frac{1}{2}}(\mathbf{Z} - \mathbf{BX})^\top$. Including the noise variance term σ^2 , the loss function for the empirical estimator is

$$\mathcal{L}_{\text{PLS}}(\mathbf{Z}, \theta) = \|\mathbf{S} - \mathbf{C} - \sigma^2 \mathbf{I}_n\|_2^2 + \operatorname{tr}[(\mathbf{Z} - \mathbf{BX})\mathbf{K}^{-1}(\mathbf{Z} - \mathbf{BX})^\top]. \quad (15)$$

208 Minimizing \mathcal{L}_{PLS} w.r.t. \mathbf{Z} and θ , we derive the empirical PLS estimators $\hat{\mathbf{Z}}_{\text{PLS}}$ and θ_{PLS} . Here, the least square
 209 term estimates the embedding \mathbf{Z} using the sample covariance, while the second term regularizes \mathbf{Z} with
 210 the kernel matrix \mathbf{K} and mean values \mathbf{BX} constructed from \mathbf{X} .

211 4.2 Maximum a posteriori (MAP) estimation

212 Another common approach to estimating covariance matrices uses maximum likelihood methods. We
 213 have already defined the data distribution \mathbf{Y} given the latent embedding \mathbf{Z} in eq. 12 and described the
 214 prior over \mathbf{Z} in eq. 8.

215 The joint distribution is then

$$p(\mathbf{Y}, \mathbf{Z} | \mathbf{X}, \theta) = \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{Z}, \rho, \sigma^2) \prod_{j=1}^d p(\mathbf{z}_j | \mathbf{X}, \mathbf{B}, r, \delta) \quad (16)$$

$$= \prod_{t=1}^T \mathcal{N}(\mathbf{y}_t | 0, \mathbf{C} + \sigma^2 \mathbf{I}_n) \prod_{j=1}^d \mathcal{N}(\mathbf{z}_j | \mathbf{b}_j \mathbf{X}, \mathbf{K}), \quad (17)$$

216 where \mathbf{C} is a function of \mathbf{Z} and ρ . Thus, the loss function for the maximum a posteriori (MAP) estimator is

$$\begin{aligned} \mathcal{L}_{\text{MAP}}(\mathbf{Z}, \theta) &= -\log p(\mathbf{Y}, \mathbf{Z} | \mathbf{X}, \theta) \\ &= \operatorname{tr}[(\mathbf{C} + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{S}] + T \log |\mathbf{C} + \sigma^2 \mathbf{I}_n| + \operatorname{tr}[(\mathbf{Z} - \mathbf{BX})\mathbf{K}^{-1}(\mathbf{Z} - \mathbf{BX})^\top]. \end{aligned} \quad (18)$$

217 Minimizing \mathcal{L}_{MAP} w.r.t. \mathbf{Z} and θ , we derive the MAP estimators $\hat{\mathbf{Z}}_{\text{MAP}}$ and θ_{MAP} .

218 4.3 Exact inference by block coordinate descent

219 Calculating the log posterior (eq. 18) for MAP inference has a computational complexity of $O(n^3)$, due
220 to the need to compute the inverse and the determinant of the covariance matrix. This cost is often
221 impractical in fMRI settings, where the number of voxels n may be on the order of thousands to hundreds
222 of thousands.

223 To optimize \mathbf{Z} and the model hyperparameters, we could use gradient descent or Newton's method
224 as the optimizer; however, this approach is computationally impractical. Thus, we need to consider a
225 scalable inference method. Existing scalable inference methods for large datasets [17–19] exploit low-rank
226 approximations to the full Gaussian process. However, these approximations suffer from a loss of accuracy
227 in covariance estimation. Thus, we develop a block coordinate descent (BCD) algorithm as an exact
228 inference method for the brain kernel model (see Methods, Algorithm 1). Coordinate descent has been
229 successfully applied to solve penalized regression models [20], to estimate covariance graphical lasso
230 models [21], and to compute large-scale sparse inverse covariance matrices [4]. Our PLS and MAP
231 estimators are non-convex smooth functions. We apply an iterative block coordinate descent method
232 solved by the proximal Newton approach [22]. Given such a scalable optimizer, we alternate between
233 optimizing \mathbf{Z} and the hyperparameters using either eq. 15 (PLS) or eq. 18 (MAP) as the objective loss
234 function. More details about the optimization can be found in the Methods section.

235 In practice, we used the PLS estimate to initialize the MAP estimate, as PLS optimization is faster and
236 requires less memory, but the MAP estimate is more principled and achieves a higher accuracy. We used
237 the BCD algorithm to optimize both the PLS and MAP objectives.

238 4.4 Predicting activity for new voxels

Although we fit the brain kernel model to data collected with a particular grid of voxels, our framework allows us to apply the model to fMRI measurements collected using different voxel grids. To do so, we use the fact that the brain kernel is defined using a Gaussian process; the mean of this GP provides a smooth mapping from 3D voxel space to the d -dimensional latent embedding space, which can be evaluated at any 3D brain locations. We obtain the embedding location of a novel 3D voxel location \mathbf{x}^* using the posterior mean of this GP:

$$\mathbf{z}^* = \mathbf{B}\mathbf{x}^* + (\mathbf{Z} - \mathbf{B}\mathbf{X})\mathbf{K}^{-1}\mathbf{k}^*, \quad (19)$$

239 where $\mathbf{k}^* = [k_f(\mathbf{x}^*, \mathbf{x}_1), \dots, k_f(\mathbf{x}^*, \mathbf{x}_n)]^\top$ represents the vector formed by evaluating the RBF covariance
240 function for the voxel at location \mathbf{x}^* and all the observed voxels in \mathbf{X} . The brain kernel for any arbitrary
241 pair of voxel locations in the test set \mathbf{x}_i^* and \mathbf{x}_j^* is therefore given by:

$$\kappa_{BK}(\mathbf{x}_i^*, \mathbf{x}_j^*) = \rho \exp(-\frac{1}{2} \|\mathbf{z}_i^* - \mathbf{z}_j^*\|_2^2), \quad (20)$$

242 with \mathbf{z}_i^* and \mathbf{z}_j^* the corresponding latent embeddings (eq. 19).

243 4.5 Synthetic experiments

244 To illustrate the performance of our proposed inference method to optimize the brain kernel objective
245 function, we began with an application to simulated data, where the ground-truth embedding is known.
246 We first compared our block coordinate descent (BCD) method (Algorithm 1), which maximizes the exact
247 log evidence (eq. 18), with two variational inference methods that optimize a lower bound on log evidence:
248 an inducing-point method [17], and a dynamical variational inference method (dynamic VIP¹) [18].

¹<https://github.com/SheffieldML/GPmat>

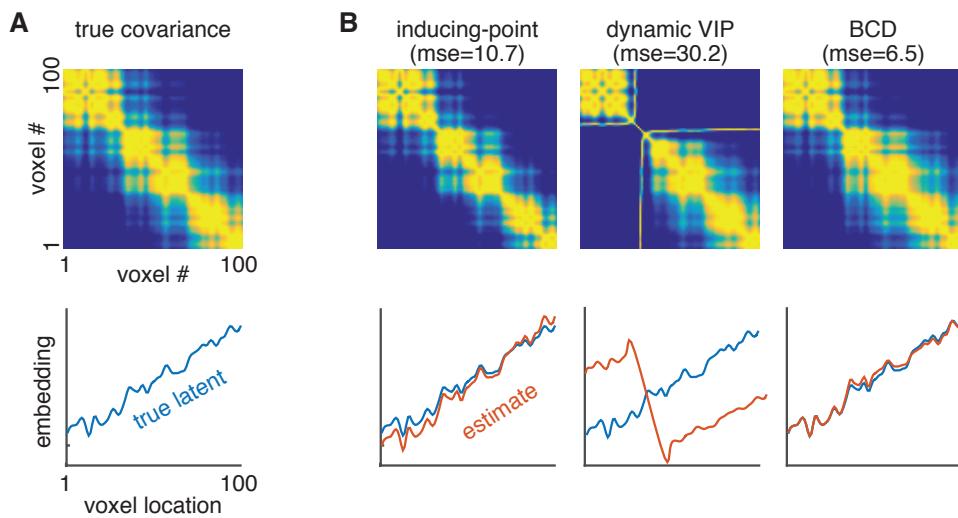


Fig 2. Recovery of 1D brain kernel from synthetic data. **A.** The true covariance matrix for neural activity at 100 evenly-spaced voxels in a 1D brain (top), generated by a 1D latent embedding function sampled from the brain kernel model (bottom, blue curve). **B.** Model estimates. The first column is the inducing-point method; the second column is dynamic VIP; and the last column is our BCD method. In each column, we show the estimated covariance matrix (top) and the estimated 1D latent embeddings (bottom, red curve). We also show the mean squared error (MSE) between the estimated covariance matrix and the true covariance matrix.

We created a simulated dataset using a 1-dimensional brain and 1-dimensional latent embedding function, sampled from the brain kernel model (Fig. 2). Here, the latent embedding is simply a nonlinearly warped version of the 1D brain. We considered a set of 100 voxel locations on an evenly spaced 1D grid: $\mathbf{X} = [1, 2, \dots, 100]^\top$. We then sampled the voxels' latent locations \mathbf{Z} from a GP with mean $m(\mathbf{x}) = 0.6\mathbf{x}$ and RBF covariance with length scale $\delta = 10$ and marginal variance $r = 9$: $\mathbf{Z} \sim \mathcal{N}(0.6\mathbf{X}, \mathbf{K})$, where $\mathbf{K}_{ij} = 9 \exp(-{(i-j)^2}/{200})$. Given this embedding function, the brain kernel defines a covariance matrix for neural activity at these 100 voxels, denoted \mathbf{C} , with the i, j th entry given by $C_{ij} = \exp(-{(\mathbf{z}_i - \mathbf{z}_j)^2}/{2})$ (Fig. 2A top). To obtain simulated fMRI measurements, we sampled 750 observations from a Gaussian distribution with zero mean and covariance $\mathbf{C} + 5I$, where $\sigma^2 = 5$ represents the variance of additive measurement noise.

We compared the different estimators on the task of recovering the true covariance and latent embedding function from this dataset (Fig. 2B). The inducing-point method with six inducing points (column 1) performed well at recovering the latent embedding function, although it was outperformed by our BCD estimator in terms of mean squared error (BCD; column 3). The dynamic VIP estimate with six inducing points (column 2) converged to a local optimum far from the true latent embedding, yielding a substantially higher error. In contrast, exact inference using BCD outperformed both inducing point-based approximate methods in terms of accuracy at recovering the true brain kernel from simulated data.

Next, we conducted a set of synthetic experiments to examine how well different models captured the covariance of simulated fMRI data, using voxels on a 1-dimensional, 2-dimensional, or 3-dimensional grid. We fit these simulated datasets using the brain kernel model optimized with (1) BCD, (2) variational inducing-point, and (3) the dynamic VIP method, and two additional models: (4) a linear brain kernel (LBK) model, and (5) a GP with RBF covariance function (RBF). Note that the first three methods are based on the original brain kernel model but have different inference methods, while the last two methods

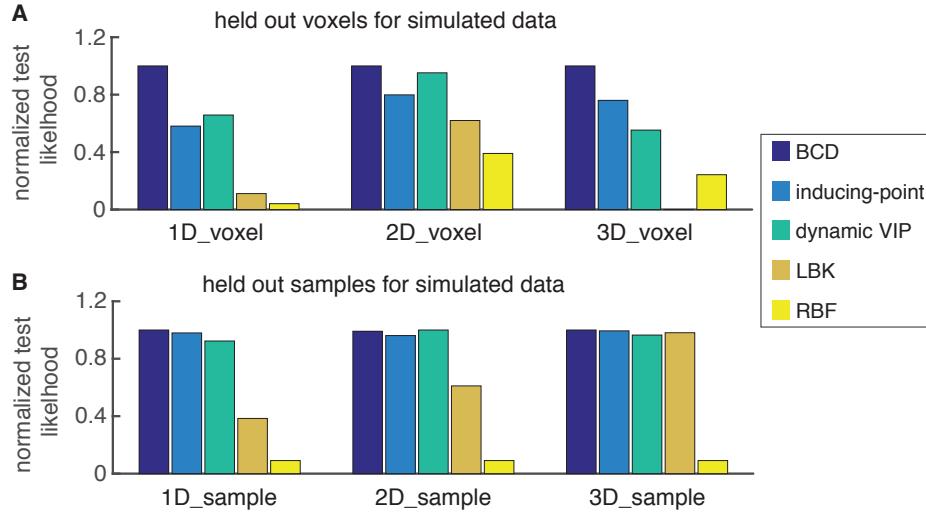


Fig 3. Quantitative comparisons for BCD, inducing-point, dynamic VIP, LBK, and RBF on three simulated datasets with 1D, 2D, and 3D input voxel locations. **A** is the held-out voxel experiment, and **B** is the held-out sample experiment. Within each experiment, we show the normalized test likelihood values for each method with three different simulated datasets.

represent simplifications of the proposed brain kernel model. The linear brain kernel model assumes a purely linear embedding function; it thus allows rotating and linearly dilating the voxel grid, but does not allow for nonlinear warping of voxel locations. The covariance function for neural activity v under the LBK model is given by:

$$\kappa_{LBK}(\mathbf{x}_i, \mathbf{x}_j) = \rho \exp\left(-\frac{1}{2}\|\mathbf{B}\mathbf{x}_i - \mathbf{B}\mathbf{x}_j\|_2^2\right), \quad (21)$$

where \mathbf{B} is a $d \times v$ linear embedding matrix and v is the number of dimensions of x , e.g., $v = 2$ for a 2-dimensional grid of voxels, and d is the dimension of the latent space. To assess the importance of the structured latent embedding in covariance estimation, we fit the data with a standard GP with RBF covariance function:

$$\kappa_{RBF}(\mathbf{x}_i, \mathbf{x}_j) = \rho \exp\left(-\frac{1}{2}\|\mathbf{x}_i - \mathbf{x}_j\|_2^2/l^2\right), \quad (22)$$

where l is the length scale for the RBF kernel. This model imposes smoothness using Euclidean distance, without estimating any (linear or nonlinear) transformation of the true voxel locations. To fit the LBK and RBF models, we maximized the following negative log likelihood for hyperparameters θ :

$$\begin{aligned} \mathcal{L}(\theta) &= -\log p(\mathbf{Y}|\mathbf{X}, \theta) \\ &= \text{tr}[(\mathbf{C}(\mathbf{X}, \theta) + \sigma^2 \mathbf{I}_n)^{-1} \mathbf{S}] + T \log |\mathbf{C}(\mathbf{X}, \theta) + \sigma^2 \mathbf{I}_n|, \end{aligned} \quad (23)$$

where $\mathbf{C}(\mathbf{X}, \theta)$ is the model-generated covariance matrix (given by eq. 21 for the linear brain kernel with $\theta = \{\mathbf{B}, \rho\}$, and by eq. 22 for the RBF model with $\theta = \{l, \rho\}$), \mathbf{S} is the sample covariance (eq. 14), σ^2 is the noise variance, and T is the number of samples in the simulated dataset.

For each grid dimension, we simulated ten independent datasets, each with different nonlinear embedding functions sampled from the brain kernel model. For the 1D experiments, we generated each dataset with 500 voxels and 750 samples, and embedded the 1D voxel space to a 1D latent embedding space. For the 2D experiments, we used a 25×25 voxel grid, embedded nonlinearly in a 3D latent embedding space, and generated datasets of 1000 samples, where each sample is a vector of 625 noisy measurements of brain activity at voxel locations. For the 3D datasets, we used a $10 \times 10 \times 10$ grid of voxels, embedded nonlinearly into a six-dimensional latent space, and we generated 1500 samples per dataset.

293 To compare models, we performed two different cross-validation tests: (i) prediction on held-out voxels
294 (Fig. 3A), and (ii) predictions on held-out samples (Fig. 3B).

295 For the first of these, we removed ten randomly-selected voxels \mathbf{X}^* from the simulated 1D datasets and
296 set them aside as test data, and we optimized the model parameters on a training set of measurements
297 from the remaining 490 voxels. (For experiments with 2D and 3D grids, we set aside 62 and 100 voxels as
298 held-out data, and we trained models using the remaining 563 and 900 voxels, respectively.) For the BK
299 and LBK models, we computed the embedding locations \mathbf{Z}^* for the test voxels using the GP posterior
300 mean given the inferred embedding locations \mathbf{Z} (eq. 19), and we used these locations to evaluate the
301 covariance of the test voxel activity (eq. 12). (For the RBF model, there is no embedding, so the covariance
302 of the test data depends only on the test voxel locations \mathbf{X}^* .) We used the resulting predictive covariances
303 to compute the log likelihood of the test data. We found that the BCD-optimized brain kernel estimate
304 outperformed other methods, while the LBK and RBF models performed worse, presumably due to their
305 inability to capture the nonlinear embedding of the simulated data (Fig. 3A).

306 Next, we evaluated cross-validation performance on held-out samples, which were measured at the same
307 set of voxels as the training set. For these simulations, we randomly selected 75 samples as test data
308 for the 1D datasets, and used the remaining 675 samples to fit the five models. (For 2D and 3D grids,
309 we used a train-test split of 900:100 and 1350:150 samples, respectively.) We estimated a covariance
310 function using measurements in the training set, and computed a test likelihood using this covariance
311 function on the test set. We generated ten random splits for each dataset and computed the normalized
312 test log likelihood. In this predictive task, the BCD estimate again outperformed other methods (Fig. 3B).

313 Overall, these simulations demonstrated that the inducing point-based GP methods, which are often used
314 due to their scalability, had higher error than our exact BCD inference method. In addition, by comparing
315 to the linear brain kernel and the standard GP model with an RBF kernel, we showed that the ability to
316 capture a nonlinear transformation of the voxel locations is critical for accurately modeling the covariance
317 of simulated brain activity.

318 5 Inferring the brain kernel from resting-state fMRI data

319 Now that we have described the brain kernel model and validated our inference method using simulated
320 data, we turn to the problem of inferring the brain kernel from real data. We fit the brain kernel model to
321 large-scale publicly-available resting-state fMRI data from the Human Connectome Project (HCP) [6]. An
322 advantage of this approach is the large size of the HCP sample, which mitigates overfitting and allows
323 us to learn an embedding function that captures correlation patterns common to a vast collection of
324 different brains. However, applying the resulting brain kernel to task-based fMRI datasets assumes that
325 the correlations presented in resting-state fMRI data are applicable to activation patterns in other brain
326 states. Previous studies have shown interesting relations between brain activity during a task and at rest.
327 One study showed that coherent spontaneous activity accounted for variability in event-related BOLD
328 responses [23]. A second study concluded that functional networks used by the brain in action were
329 continuously and dynamically “active” even when at “rest” [24]. A third study showed that resting-state
330 activation patterns had strong statistical similarities to cognitive task activation patterns [25]. These
331 provide reasons for optimism, though the possibility of changes in correlation across different tasks could
332 potentially affect the application of the brain kernel which we will show empirically later.

333 We examined resting-state fMRI data from 812 subjects collected via the Human Connectome Project
334 (HCP) [6]. These data were acquired in four fMRI runs of approximately 15 minutes each, two runs in one
335 session and two in another session, with eyes open and relaxed fixation on a projected bright cross-hair
336 on a dark background. A sophisticated preprocessing pipeline was used to align voxels across subjects.
337 Detailed information about imaging protocols, image acquisition, and preprocessing can be found in [26].

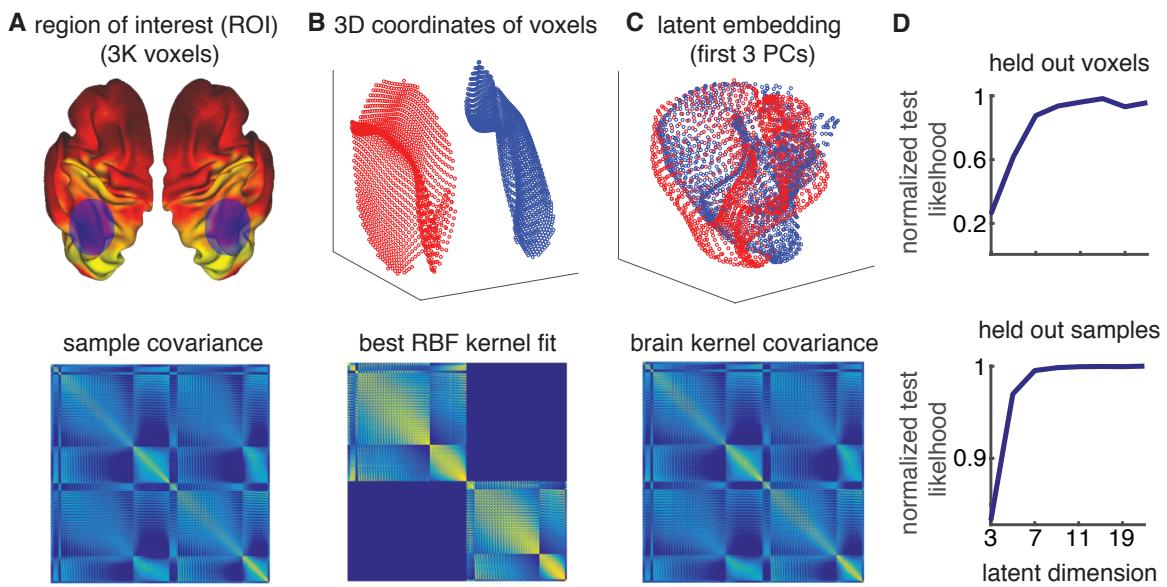


Fig 4. 3D embeddings and the corresponding covariance matrices of resting-state fMRI data. **A.** The two example regions of interest (ROIs) we selected in the left and right parietal lobes (top) and the full covariance for the 3K voxels in these two ROIs (bottom). **B.** The original 3D voxel coordinates (top) and the best-fitting RBF kernel (bottom). **C.** A 3-dimensional projection of the estimated 20-dimensional brain kernel embedding (top) and the corresponding brain kernel covariance matrix (bottom). **D** shows the influence of the latent dimension on the predictive performance for held-out voxels and held-out samples with the resting-state fMRI data.

338 The resulting dataset consisted of 59,412 voxels in a cortical surface coordinate system. Although the full
339 covariance matrix of these data is of size $\approx 59K \times 59K$, we fit the model using the first 4500 eigenvectors
340 of the full matrix provided by HCP.

341 We fit the brain kernel with different numbers of latent dimensions and computed the test log likelihood
342 with held-out voxels and samples (Fig. 4D). We found that test performance plateaued with increasing
343 dimensionality, and selected $d = 20$ dimensions for subsequent analyses. We then estimated the brain
344 kernel by fitting the 20-dimensional latent embedding for each voxel using BCD optimization of the log
345 marginal likelihood. The resulting function is a matrix of embedding locations of size $\approx 59K \times 20$,
346 where each row contains the embedded location of a single voxel in the resting-state fMRI dataset. This
347 embedding, in addition to the hyperparameters (the linear projection matrix \mathbf{B} and the hyperparameters
348 $\{\gamma, \delta\}$ for \mathbf{K}), provide a full parametrization of the brain kernel.

349 Although it is impossible to visualize the 20-dimensional nonlinear embedding that defines the brain
350 kernel, we can gain insight into its shape by plotting low-dimensional projections on subsets of voxels. To
351 visualize the brain kernel, we selected two symmetric ROIs in the left and right parietal lobes (Fig. 4A,
352 top). Taken together, these two ROIs contain 3K voxels. We visualize the sample covariance of these
353 voxels, computed directly from the eigenvectors (Figure 4A, bottom). This covariance contains four
354 identifiable blocks: two diagonal blocks that correspond to the covariance of voxels within each ROI, and
355 two off-diagonal blocks that correspond to cross-ROI covariances. The off-diagonal blocks reveal that the
356 two ROIs have reasonably strong correlations despite being spatially distant in the brain.

357 The original 3D voxel coordinates (Fig. 4B, top) and the covariance given by the best-fitting RBF kernel
358 (eq. 22; Fig. 4B, bottom) show that the RBF kernel model fails to capture the the off-diagonal blocks of the

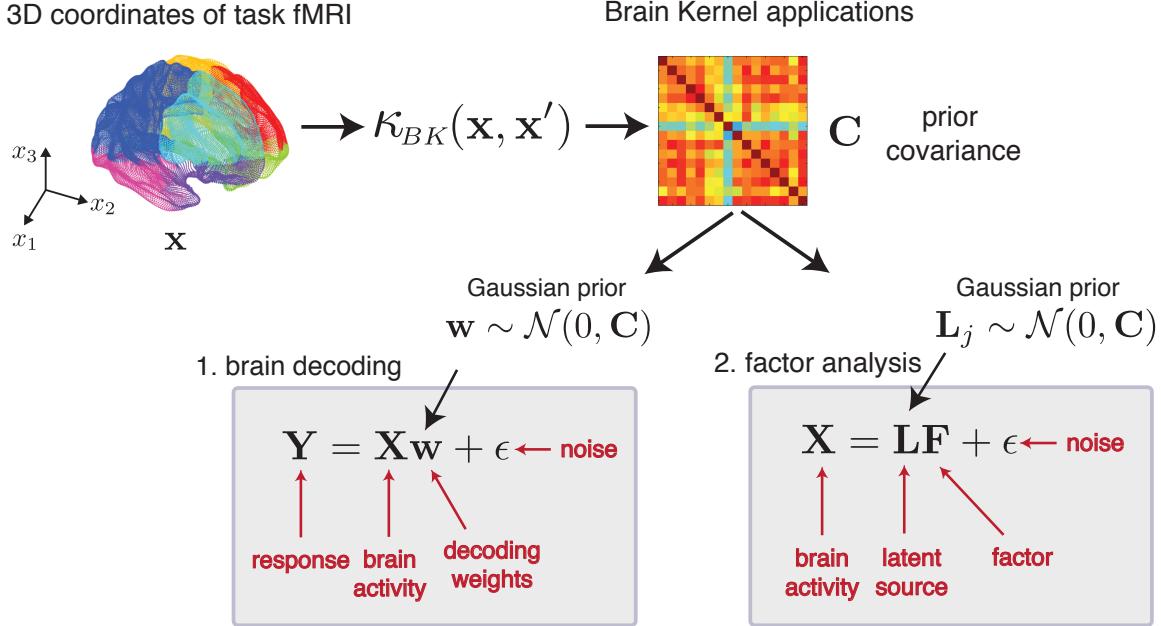


Fig 5. Schematic figure illustrating two applications of the brain kernel to task fMRI data. After fitting the brain kernel to resting-state data, we applied it to task fMRI data with n voxels by evaluating the brain kernel at the 3D voxel locations. This results in a $n \times n$ prior covariance matrix for the task data, denoted \mathbf{C} . We then used the covariance \mathbf{C} as the prior covariance for two modeling tasks: (1) brain decoding and (2) factor analysis. Two unknown parameters \mathbf{w} and \mathbf{L} are both random variables with a Gaussian prior whose covariance is \mathbf{C} . Thus, we effectively imposed assumptions on the structure of \mathbf{w} and \mathbf{L} via \mathbf{C} .

359 sample covariance, corresponding to covariances between voxels in opposite hemispheres, due to the fact
 360 that the RBF kernel depends only on the Euclidean distance between voxels. In contrast, a 3-dimensional
 361 projection of the estimated 20-dimensional brain kernel embedding (Fig. 4C, top) and the corresponding
 362 brain kernel covariance (Fig. 4C, bottom) appears to capture this cross-ROI structure in the covariance
 363 matrix. The 3D projection corresponds to the first three principal components of the 20-dimensional latent
 364 embeddings, and shows that paired voxels from opposite hemispheres are embedded close to each
 365 other under the brain kernel. Conversely, some voxels that are physically close together in the brain are
 366 embedded far apart in the embedding space. This nonlinear embedding allows the brain kernel to more
 367 accurately capture the covariance of the real data.

6 Applications

369 We applied the brain kernel to two different fMRI data analyses (Figure 5). Note that we estimated the
 370 brain kernel using the resting-state fMRI from the HCP database as described above. We kept the same
 371 brain kernel across all analyses. We didn't re-estimate the brain kernel using any task fMRI. For both
 372 applications, we took the 3D brain coordinates for n voxels from a task fMRI dataset of interest. We
 373 input these coordinates and corresponding observations into the brain kernel, which produced a $n \times n$
 374 covariance matrix \mathbf{C} (using eq. 19 and eq. 20). We then used this covariance as the prior covariance for
 375 Bayesian analyses for different modeling tasks, brain decoding and factor analysis. We describe these
 376 applications in detail below.

377 6.1 Brain decoding

378 In this section, we illustrate how the brain kernel can be applied to fMRI classification (or "decoding")
379 tasks.

380 6.1.1 HCP tasks

381 We first examined the task fMRI datasets in the HCP database. We explored the working memory task,
382 the gambling task [27], the language processing task [28], the motor task [29], the emotion processing
383 task [30], the relational processing task [31], and the social cognition task (more details found in [26, 32]).
384 We will elaborate on the working memory and gambling tasks and finally summarize the result with all
385 datasets.

386 Working memory task

387 We obtained the working memory task fMRI measurements from the HCP [32]. The stimuli consisted of
388 four types of pictures: places, tools, faces, and body parts. Stimuli were presented for 2 s on each trial
389 followed by a 500-ms inter-trial interval (ITI). Each task block consisted of ten trials of 2.5 s each. Each
390 run contained 8 task blocks, half for a 2-back working memory task and half for a 0-back working memory
391 task, along with 4 fixation blocks. Two runs were collected for each subject, with 405 volumes per run or
392 approximately 5 minutes. Because the task fMRI was from the same HCP project as the resting-state
393 fMRI, they shared the same coordinate system and preprocessing pipeline. The task fMRI was aligned
394 with the MNI template with the same 59,412 voxels as the resting-state fMRI data. Instead of analyzing the
395 whole-brain data for brain decoding, we worked with functional ROIs. For each presented type of pictures,
396 we took the group-average task contrast for the 2-back vs 0-back working memory task and kept all voxel
397 coordinates whose z-statistics were -1.96 or +1.96 standard deviations away from the mean, indicating
398 that the voxels were statistically significant in the contrast map. We repeated the same thresholding
399 procedure for all objects and took a union set of all coordinates to formulate the functional ROIs for the
400 working memory task. We didn't select the ROIs using the contrast among objects, therefore the resulting
401 ROIs contained no discriminative information with respect to the decoding task.

402 The experiment required observers to perform a working memory task using four different types of objects.
403 We tested the ability to decode these objects from fMRI data by fitting a binary linear classifier for each
404 pair of objects. We used Bayesian linear regression classifiers to solve these six binary classification
405 problems. A Bayesian linear regression classifier has the form,

$$y = \text{sign}(\mathbf{x}\mathbf{w} + \epsilon), \quad (24)$$

where \mathbf{x} is a vector of all voxels for one fMRI measurement or sample, y is a ± 1 label indicating the binary object category for that sample, \mathbf{w} is a vector of regression coefficients, and ϵ is independent zero-mean Gaussian noise with variance σ^2 . To regularize the estimate of \mathbf{w} , we assumed a zero-mean Gaussian prior, $\mathbf{w} \sim \mathcal{N}(0, \mathbf{C})$, with prior covariance \mathbf{C} . We considered three different choices of prior covariance: (1) a ridge prior, which corresponds to a diagonal covariance with a positive constant along the diagonal; (2) a radial basis function covariance (eq. 22), which imposes smoothing based on the voxels' 3D locations in the brain, and (3) a covariance derived from the brain kernel. For the brain kernel prior, we used the covariance function

$$\kappa_{BK}(\mathbf{x}_i, \mathbf{x}_j) = \rho \exp\left(-\frac{1}{2}\|\mathbf{f}(\mathbf{x}_i) - \mathbf{f}(\mathbf{x}_j)\|_2^2/l^2\right), \quad (25)$$

406 where \mathbf{f} is the nonlinear embedding function from the brain kernel, and $\{\rho, l\}$ are tuneable hyperparameters.

408 We trained the classifier with these three priors on one run and calculated accuracy performance on the
409 second run, then repeated the same procedure with test and training sets reversed. When we trained

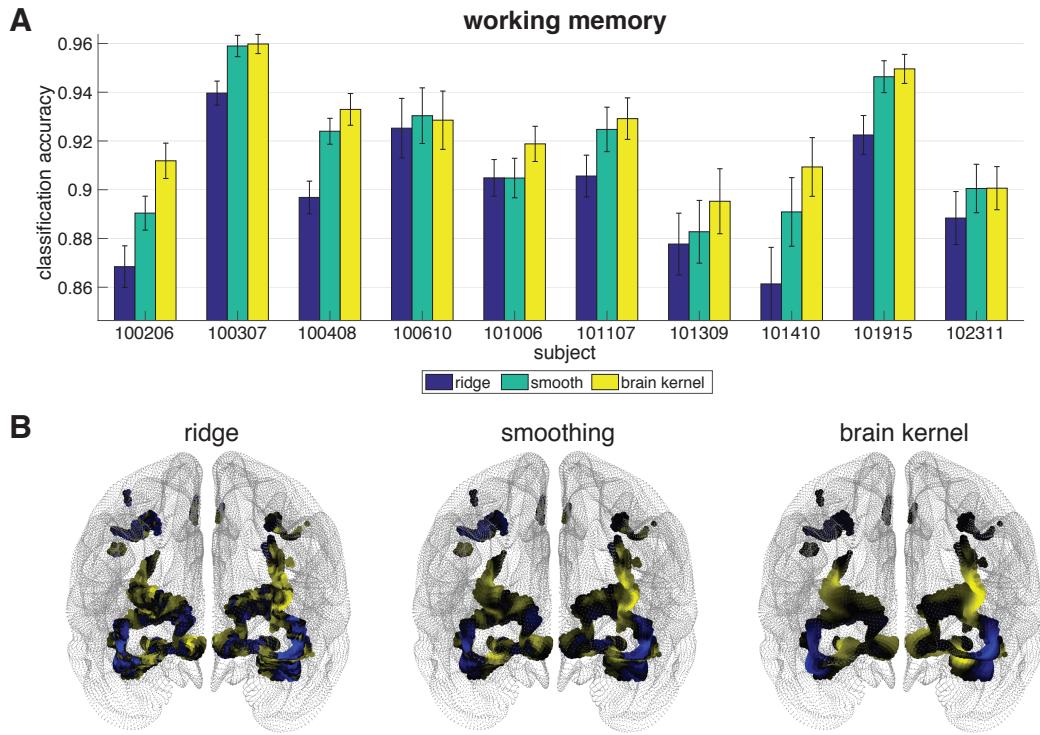


Fig 6. A. Accuracy performance on the working memory task. The x-axis indicates subject identifiers. The y-axis is accuracy performance. We compared our brain kernel with a ridge prior and a smooth RBF kernel, color coded. The error bars indicate standard errors. **B.** Visualization of an example set of decoding weights. Blue indicates negative values and yellow indicates positive values.

the model, we randomly split the training run into five folds and selected the optimal hyperparameters in the covariance functions via 5-fold cross validation. After cross validation, we applied the optimal hyperparameters to the test run for each prior model. We repeated this 5-fold cross validation experiment ten times to reduce variability. We used linear regression to train the classifier, so the ± 1 labels were treated as continuous target values, and test accuracy was evaluated by taking the sign of the prediction. We plotted averaged accuracy across both runs and all repeats for the three priors for ten randomly-selected subjects (Fig. 6A). The RBF prior outperformed the ridge prior, but the brain kernel prior outperformed both of the other priors, indicating that smoothing in a nonlinear embedding space defined by correlations of fMRI signals provided additional benefits in regularizing weights for a classification task. Moreover, the framework of the brain kernel for regularization allows us to visualize the inferred decoding weights overlaid on a 3D brain (Figure 6B).

421 **Gambling task**

422 We next examined fMRI data in a gambling task from the HCP [32], adapted from a prior study [27].
 423 Participants were asked to play a card-guessing game. They were shown a mystery card with a number
 424 that could range from 1 to 9. They needed to guess whether it was more or less than 5 by pressing on of
 425 two buttons. If they made the correct guess, the card showed a green up arrow with "\$1" for rewards; if
 426 they guessed wrongly, the card showed a red down arrow with "-\$0.5" for losses; if the true value was 5,
 427 they got a neutral response without win or loss. Participants had 1500 ms to guess, and the feedback
 428 was presented for 1000 ms, followed by a 1000-ms inter-trial interval (ITI). Each task block consisted
 429 of eight trials that were either mostly reward or mostly loss. Two runs were collected for each subject.

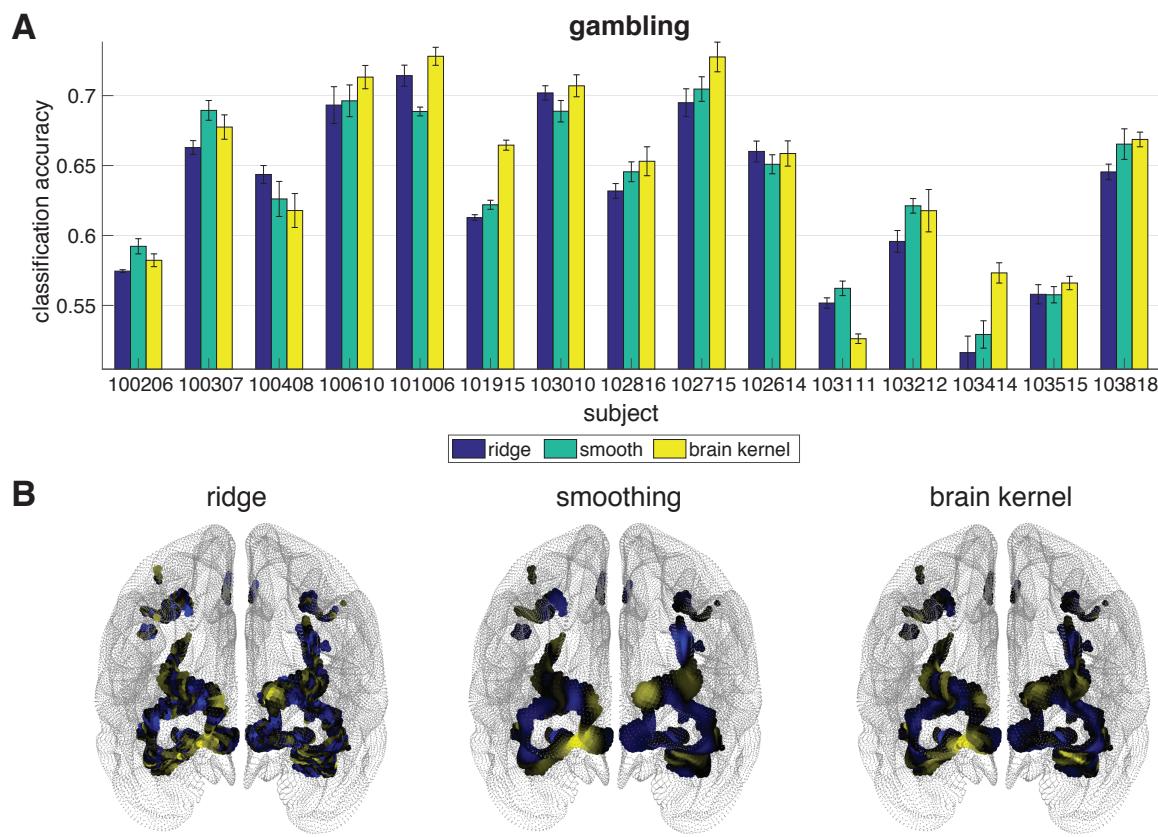


Fig 7. A. Accuracy performance on the gambling task. The x-axis indicates subject identifiers. The y-axis is accuracy performance. We compared our brain kernel with a ridge prior and a smooth RBF kernel, color coded. **B.** Visualization of an example set of decoding weights. Blue indicates negative values and yellow indicates positive values.

430 Each run contained two mostly reward and two mostly loss blocks, interleaved with four fixation blocks.
431 There were 253 volumes per run, which lasted approximately 3 minutes. The task fMRI was aligned to the
432 MNI template and had the same 59,412 voxels as the resting-state fMRI data. Instead of analyzing the
433 whole-brain data for brain decoding, we worked with functional ROIs which were selected using the same
434 approach as described in the working memory task.

We formulated the task as a binary classification problem separating reward trials and punishment trials. We used the same Bayesian linear regression classifiers as described in the working memory section. We trained the classifier with three priors on one run and calculated the accuracy performance on the second run, then switched the training and test runs. This procedure was repeated ten times. The ± 1 labels were treated as continuous target values during training, and the test accuracy was evaluated by taking the sign of the prediction. We computed the averaged accuracies across two runs and 10 repetitions for the three priors for 15 subjects (Fig. 7A). For 11 out of 15 subjects, the brain kernel improved accuracy over both the ridge prior and the smooth RBF prior. For two of the remaining four subjects, the brain kernel outperformed the ridge prior. The discrimination problem with the gambling data was more difficult than the working memory task. Many activations occur in the visual cortex and the prefrontal cortex for working memory, whereas critical activations for the reward task may be localized to the striatum, which is not included in the cortical data used here. However, with the cortical brain kernel, we were still able to improve the predictive ability of the Bayesian model.

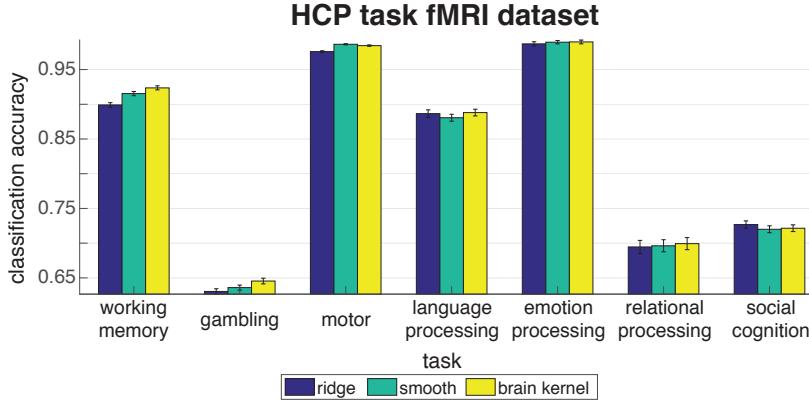


Fig 8. Averaged accuracy performance for all task fMRI datasets in the HCP database. The x-axis indicates the task. The y-axis is accuracy performance.

448 All HCP tasks

449 We've elaborated on the working memory and gambling tasks above. We also achieved the classification
450 accuracy performance for all other tasks (presented in Appendix B). Here we summarize the averaged
451 accuracy over all subjects for each task in Fig. 8. Consistent with the above results, we succeeded in
452 achieving the best performance with the working memory and gambling tasks using the brain kernel. For
453 other tasks, the brain kernel performed mildly better than the ridge prior and the smoothing prior estimates.
454 Regularizing the decoding weights with the brain kernel did not improve performance in a dominating way,
455 suggesting that the covariance of resting-state fMRI did not provide significantly useful information for
456 classifying fMRI measurements in these tasks, or at least that the brain kernel model was not capable of
457 fully capturing it. It remains possible, however, that a brain kernel trained on task fMRI datasets might
458 offer benefits for decoding fMRI data from these tasks.

459 6.1.2 Visual recognition task

460 Next, we examined the problem of decoding faces and objects from fMRI measurements during a visual
461 recognition task. Just to remind, the brain kernel was estimated using the resting-state fMRI from HCP,
462 and we applied it to a popular fMRI dataset from a study of human ventral temporal cortex [33] for the
463 decoding task. We extended the application beyond HCP, i.e., constructing the brain kernel on HCP and
464 then trying on a completely different dataset. Therefore, we were looking at across-dataset generalization,
465 not just across-task generalization within HCP. In this visual recognition experiment, six subjects were
466 asked to recognize eight different types of objects (bottles, houses, cats, scissors, chairs, faces, shoes,
467 and scrambled control images, examples in Fig. 9). Each subject participated 12 scanning runs. In
468 each run, the subjects viewed images of eight object categories, with 11 whole-brain measurements per
469 category. Each subject's fMRI data was preprocessed using the fMRIprep package² [34] and aligned to
470 the MNI template. Both voxels in this dataset and voxels in the HCP database were all aligned in the same
471 MNI space, allowing us to use eq. 19 to get the brain kernel covariance for the present dataset. Instead of
472 analyzing the whole-brain data, we extracted ROIs with 1645 voxels in the ventral temporal cortex, which
473 is thought to be involved in object recognition. The ROI mask was obtained from Nilearn [35].

474 We assessed performance by training Bayesian linear regression classifiers to discriminate between pairs

²<https://github.com/poldracklab/fmriprep>

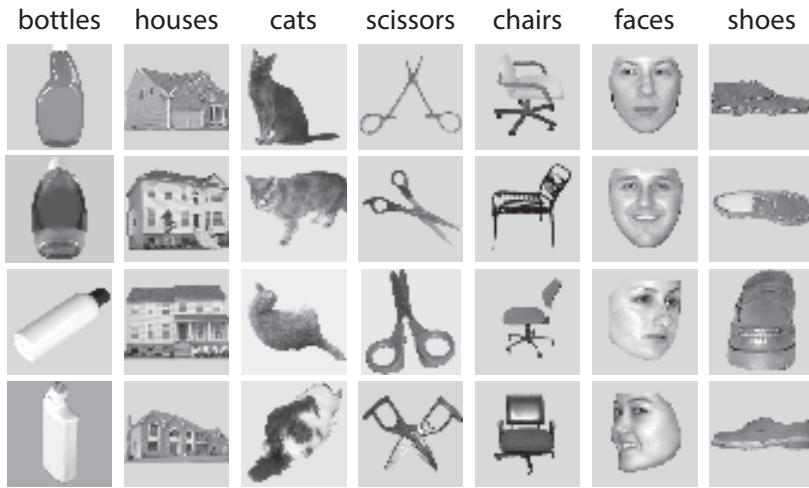


Fig 9. Examples of the stimuli for 7 categories (except for scrambled control images) [33].

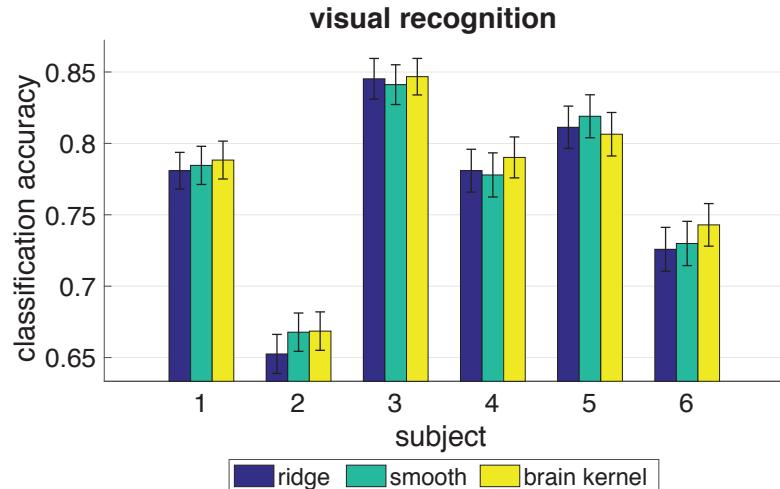


Fig 10. Accuracy performance on the visual recognition task. The x-axis indicates subject IDs. The y-axis is accuracy performance. We compared our brain kernel with a ridge prior and a smooth RBF kernel, color coded.

of objects, e.g., face vs. bottle, for each of the 28 possible binary classifications among the eight objects (Fig. 9). We trained the weights w for each model using linear regression from fMRI measurements x to binary labels $y \in \{-1, +1\}$, and assessed accuracy on the test set using predicted labels $\hat{y} = \text{sign}(xw)$. Here, we split the training and test sets by subjects. In this visual recognition dataset, we had six subjects but only $11 \times 2 \times 12 = 264$ measurements in a 1645-voxel space for each subject in a binary decoding task. The number of training measurements was not sufficiently large to train a good classifier and lead to a reasonable generalization performance. Therefore, different from the HCP tasks, we chose to do inter-subject analyses by using 5 subjects for training and one subject for test. We repeated this leave-one-subject-out manner for six times with each subject being used as the test set once and obtained the result in Fig. 10.

485 The averaged accuracy performance across four repeated runs for six subjects shows that the brain kernel
486 performed comparably to the ridge and smooth priors with better accuracy performance for 5 out of 6
487 subjects (Fig. 10). This indicates that the brain kernel can provide functional and structural support for
488 most subjects and tasks. The improvement was statistically modest overall based on the standard errors,
489 which could be a result of several factors: misalignment of the coordinate space to the HCP coordinate
490 space used to estimate the brain kernel; mismatch between the resting-state covariance used to construct
491 the brain kernel and covariance present during the visual recognition task; or the object recognition tasks
492 may rely on fine-grained spatial response topographies that are poorly aligned across individuals.

493 6.2 Factor modeling

494 In this section, we illustrate a second type of application of the brain kernel. Instead of using it to regularize
495 decoding weights, as in the previous section, we use it as a spatial prior for Bayesian factor analysis.

496 6.2.1 Sherlock movie watching task

497 We examined the Sherlock fMRI dataset, in which participants were scanned while they watched the
498 British television program “Sherlock” for 50 min [36]. The fMRI data comprised 1,973 TRs (Repetition
499 Time), where each TR was 1.5 s of the movie. Before performing any analysis, the fMRI data were
500 preprocessed and aligned to MNI space using the techniques described in prior work [36]. We examined
501 the brain data averaged across all subjects to smooth out individual variability. We identified 11 ROIs
502 previously implicated in processing naturalistic stimuli, comprising the default mode network (DMN-A,
503 DMN-B), the ventral and dorsal language areas, and the primary auditory and visual cortices [37].

504 For each ROI, we performed a standard factor analysis (FA) to factorize the voxel-by-time fMRI data into a
505 latent source matrix and a factor matrix. Thus fMRI images can be considered as being generated by a
506 covariate-dependent superposition of latent sources. Some following analyses, such as decoding and
507 encoding tasks, can be performed with the factor matrix. The number of latent sources is much fewer than
508 the number of time points, thus providing a parsimonious description of neural activity patterns that avoids
509 many of the pitfalls of traditional voxel-based approaches. Similar FA based models have been proposed
510 in [38, 39]. Here, we performed a Bayesian factor analysis of the form

$$\mathbf{X} = \mathbf{LF} + \epsilon, \quad (26)$$

511 where $\mathbf{X} \in \mathbb{R}^{N \times T}$ is the fMRI image with N voxels and T time points, $\mathbf{L} \in \mathbb{R}^{N \times K}$ is the latent source
512 matrix encoding the canonical spatial pattern (over voxels) associated with each latent source. $\mathbf{F} \in \mathbb{R}^{K \times T}$
513 is the factor matrix containing the timeseries for K latent sources. ϵ is a Gaussian noise with $\epsilon \sim \mathcal{N}(0, \sigma^2)$.
514 We assumed that the i th factor (column) is from a standard normal distribution, i.e., $p(\mathbf{F}_i) = \mathcal{N}(0, \mathbf{I})$. We
515 also assumed that the j th latent source (column) is from a Gaussian prior $p(\mathbf{L}_j) = \mathcal{N}(0, \mathbf{C})$ where \mathbf{C}
516 is the covariance matrix. Like the decoding tasks, we used a ridge covariance, a smooth RBF kernel,
517 and the brain kernel in the place of the covariance matrix \mathbf{C} . With different prior covariances, the latent
518 sources $\mathbf{L}_{:,j}$ showed different characteristics imposed by these distinct regularizations.

519 Our goal was to infer the latent variables \mathbf{L} and \mathbf{F} , and the hyperparameters for \mathbf{C} and σ^2 , denoted as θ . To
520 quantify performance, we compared data explanation under the three prior covariances with the same num-
521 ber of latent sources which was set to 10. In practice, we first standardized the fMRI image \mathbf{X} , then solved
522 the optimization of the hyperparameters using the marginal distribution, and finally inferred the factor matrix
523 and the latent sources via maximum likelihood parameter estimation (details in Methods). To evaluate the
524 performance of the three priors, we used a method called “co-smoothing” [40], depicted in Fig. 11. We first
525 split the time points into two equal sets by taking the first half as one set and the latter half as the other set.

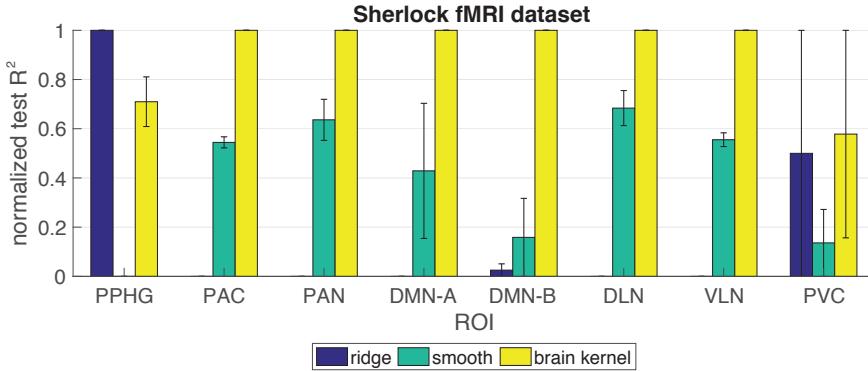


Fig 12. Normalized test R^2 performance on the Sherlock fMRI dataset. The x-axis indicates ROIs (PPHG – Posterior Parahippocampal Gyrus; PAC – Primary Auditory Cortex; PAN – Primary Auditory Network; DMN-A – Default Mode Network-A; DMN-B – Default Mode Network-B; DLN – Dorsal Language Network; VLN – Ventral Language Network; PVC – Primary Visual Cortex). The y-axis is the normalized R^2 performance on the test set (higher values indicate better performance). The error bars indicate standard errors. We compared our brain kernel with a ridge prior and a smooth RBF kernel, color coded.

We took the first set for training (the blue region) and the second set for inference (the pink region) and test (the yellow region). We trained the model with the neural activity in the training set to obtain the estimated latent sources \mathbf{L}^* . We then kept the latent sources fixed for the inference and test purpose. Next we split the second set into five folds along the voxel axis, four for inference and one for test. We used the neural activity in the inference set to infer the factor matrix (mapping the latent sources to the time series) during the inference period given the latent sources. Finally we predicted the neural activity for the left-out voxels in the test set given the latent sources and factor matrix. We repeated the inference and test step five times in a cross-validation fashion and obtained an averaged R^2 value representing the performance of the prior. After the first run, we achieved three R^2 values for the three priors (ridge, smooth and brain kernel). To better visualize the difference, we normalized the three R^2 values so that the maximum was 1 and the minimum was 0. We then launched a second run by using the second set as the training set and the first set as the inference and test set. The final normalized test R^2 value was an average of the two runs.

We compared the normalized test R^2 value for the three priors for each ROI (Fig. 12). The error bars indicate standard errors. In most regions, the brain kernel outperformed the ridge covariance and the smooth RBF kernel. This implies that when performing Bayesian FA, the brain kernel provides a superior prior covariance for the latent source matrix and may enhance performance in terms of data explanation.

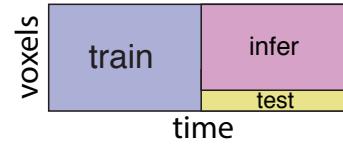


Fig 11. Co-smoothing evaluation. The blue region is the data used for training; the pink region contains voxels that are used to infer the factor matrix during the test period; and the yellow region is used for test.

6.2.2 HCP tasks

We examined the task fMRI datasets in the HCP database with the same Bayesian FA model. We collected all the task fMRIs for the same 10 subjects and performed Bayesian FA for each subject in each task. We implemented the same “co-smoothing” evaluation and compared the normalized test R^2 for the three priors averaged over all subjects for each task (Fig. 13). In most tasks, the brain kernel outperformed the

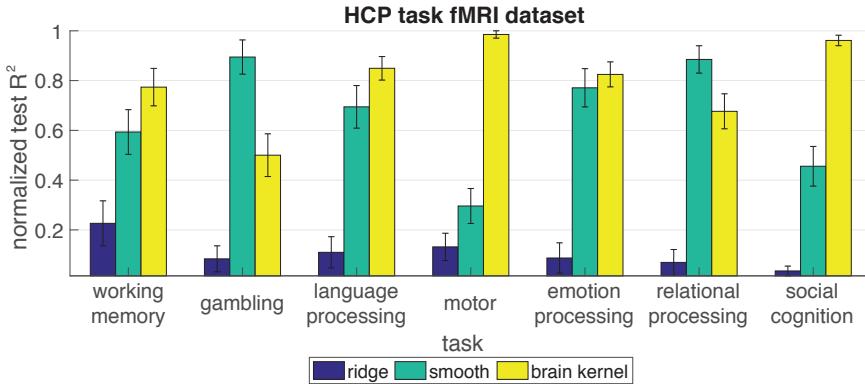


Fig 13. Normalized test R^2 performance on the HCP task fMRI datasets. The x-axis indicates the task. The y-axis is the normalized test R^2 performance averaged over 10 subjects for each task (higher values indicate better performance). The error bars indicate standard errors. We compared our brain kernel with a ridge prior and a smooth RBF kernel, color coded.

552 ridge covariance and the smooth RBF kernel. This implies that the brain kernel provides a superior prior
 553 covariance for the latent source matrix and may enhance performance in terms of data explanation for the
 554 HCP database.

555 7 Discussion

556 We introduced the brain kernel, a new model for the spatial covariance of fMRI data. This model takes the
 557 form of a covariance function, meaning that it can take any two continuous voxel locations in the brain as
 558 inputs and return the prior covariance of their activities. Unlike standard smoothness-inducing covariance
 559 functions used in the Gaussian process literature, which assume that correlation falls off monotonically
 560 with distance, the brain kernel allows widely-separated voxel locations (e.g., on opposite sides of the brain)
 561 to have high correlation, while pairs of nearby voxels can be nearly uncorrelated. This property arises
 562 from the fact that the nonlinear embedding function, parametrized with Gaussian processes, warps and
 563 stretches the brain in a higher-dimensional space so that widely-separated pairs of voxels can be mapped
 564 to nearby embedding locations, while pairs of nearby voxels can be moved far apart in embedding space.

565 The brain kernel model is a hierarchical extension of the Gaussian process latent variable model. We fit
 566 the brain kernel model using a nonlinear embedding of the 3D brain in a 20D space. We introduced an
 567 exact inference method for fitting the brain kernel using block coordinate descent (BCD), and trained it
 568 using a large publicly-available dataset of resting-state fMRI measurements from many subjects. We found
 569 that the resulting brain kernel accurately captured the covariance of resting-state fMRI measurements.

570 We further demonstrated the applicability of the brain kernel using two different modeling tasks: brain
 571 decoding and factor analysis. In both cases, we used the brain kernel to define a prior distribution in place
 572 of a more conventional prior based on ℓ_2 shrinkage or local smoothness. In both cases we showed that
 573 the brain kernel can achieve gains in performance, illustrating that the correlations in resting-state fMRI
 574 may be usefully mined to improve the accuracy of task fMRI.

575 **Limitations**

576 Despite the successes presented above, it is important to acknowledge that we applied the brain kernel to
577 several other datasets for both the decoding and the factor modeling tasks beyond the examples shown in
578 the main section, for which it did not yield improvements over standard priors.

579 For decoding, we applied the brain kernel to the Sherlock movie dataset [36] whose decoded vectors
580 contained semantic descriptions for the movie scenes. We observed that the smoothing prior and the
581 brain kernel prior both converged to the ridge prior after hyperparameter estimation, i.e., the estimated
582 length scale was very small and thus smoothness did not help the decoding task. We also applied it to a
583 public fMRI dataset in which subjects viewed 5000 visual images (BOLD5000) [41]. The binary labels were
584 living objects versus nonliving objects. We observed that the brain kernel did not reliably outperform the
585 ridge prior and the smoothing prior estimates across subjects. More details can be found in Appendix B.
586 For these cases, regularizing the decoding weights with the brain kernel did not improve performance,
587 suggesting that the covariance of resting-state fMRI did not provide useful information for classifying
588 fMRI measurements in these tasks, or at least that the brain kernel model was not capable of capturing
589 it. It remains possible, however, that a brain kernel trained on task fMRI datasets might offer benefits for
590 decoding fMRI data from these tasks.

591 For factor modeling, we also applied the brain kernel to the visual recognition task and the BOLD5000
592 dataset except for the HCP dataset and the Sherlock movie dataset presented in the paper. For most
593 subjects in these two datasets, the brain kernel didn't show a dominating performance over both the ridge
594 covariance and the smooth RBF kernel (Appendix C).

595 The lack of benefit observed on non-HCP (outside the HCP) datasets may arise from a mismatch between
596 the covariance of resting-state fMRI observations used to fit the brain kernel and the covariance of task
597 fMRI datasets. However, it might also arise from differences in acquisition parameters, preprocessing
598 steps, or alignment between HCP and non-HCP datasets. Although we aligned all voxels with the MNI
599 template, misalignment might still result from differences between processing pipelines.

600 **Outlook and future directions**

601 Although the brain kernel did not achieve improved performance in all datasets and applications we
602 considered, we feel it nevertheless holds great potential for fMRI data analysis. First, by using the
603 Gaussian process latent variable modeling (GPLVM) framework, we cast the problem of estimating
604 covariance of fMRI datasets as that of estimating a nonlinear mapping from 3D brain coordinates to a
605 latent embedding space. This results in a compact representation of the full brain covariance matrix,
606 requiring storage of only a $N_{voxels} \times 20$ matrix of embedding locations (when the embedding dimensionality
607 is 20), as opposed to a full $N_{voxels} \times N_{voxels}$ covariance matrix. Moreover, because the brain kernel is a
608 continuous covariance function, we can use it to model the covariance at arbitrary voxel locations, even
609 those not contained in the original training dataset.

610 In addition to providing a compact parametrization of fMRI covariance, the brain kernel's nonlinear
611 embedding function may be useful for gaining insights into the geometry of correlations within and across
612 brain regions. Examining the embedding of different brain regions (e.g., as shown in Fig. 4) allows
613 researchers to directly visualize correlations in terms of distances between embedded voxels.

614 Although the brain kernel fit to the HCP resting-state fMRI data did not yield improvements on all task
615 fMRI datasets we explored, it is possible that other methods for training or applying the brain kernel
616 might produce bigger gains. For example, one might train the brain kernel on task fMRI datasets, or
617 train a hierarchical version that gains statistical strength from combining datasets, while preserving
618 flexibility to capture differences between the two kinds of data. A more ambitious possibility is to formulate

619 a hierarchical version of the brain kernel that allows for per-subject variability. This would produce a
 620 hierarchical brain kernel in which each brain has own specific brain kernel, allowing for detailed differences
 621 in correlation maps across brains. All such applications will benefit from more robust methods for alignment
 622 and preprocessing, and it may be that these alone will be enough to improve the performance of the brain
 623 kernel.

624 Finally, an exciting possibility for future work is to combine the brain kernel with other advanced statistical
 625 modeling techniques. Methods for modeling structured variability, such as topographic ICA [39], and
 626 methods for structured sparsity, such as GraphNet [42], sparse overlapping sets lasso [43], and dependent
 627 relevance determination [44, 45], rely on capturing dependencies between nearby voxels. All such methods
 628 might thus be improved by using nearness in functional embedding space provided by the brain kernel,
 629 instead of nearness in 3D Euclidean space inside the brain. Likewise, methods for linear alignment of
 630 functional data from multiple subjects such as hyperalignment [46, 47] might be extended using nonlinear
 631 warping of brain coordinates under the brain kernel. Therefore, we feel that the brain kernel holds promise
 632 for inspiring new data-driven prior distributions and new modeling approaches for capturing structure in
 633 fMRI data.

634 Methods

635 Block coordinate descent for the brain kernel

636 The penalized least squares (eq. 15, PLS) has a computational complexity of $O(n^2)$ and memory storage
 637 of $O(n^2)$; however, the maximum a posteriori (eq. 18, MAP) has a computational complexity of $O(n^3)$ –
 638 to invert the covariance matrix – and memory storage of $O(n^2)$. n is the number of voxels, which often
 639 exceeds 10k. Gradient descent or Newton's method is computationally impractical as the optimizer. Thus,
 640 we need to use a scalable inference method. Existing inference methods for large datasets [17–19] exploit
 641 low-rank approximations to the full Gaussian process, which, however, suffer from a loss of accuracy
 642 in covariance estimation. Thus, in this section, we develop a block coordinate descent algorithm as an
 643 exact inference method for the brain kernel. Coordinate descent has been successfully applied to solve
 644 penalized regression models [20], to estimate covariance graphical lasso models [21], and to compute
 645 large-scale sparse inverse covariance matrices [4].

646 Our PLS and MAP estimators are non-convex smooth functions. We apply an iterative block coordinate
 647 descent method solved by the proximal Newton approach [22]. We first divide the voxel set $\{1, \dots, n\}$
 648 into blocks. Next, we iterate over all blocks, minimizing the functions with respect to the voxels within
 649 each block. Without loss of generality, we split the voxel set into two blocks, $\{1, \dots, m\}$ (block 1) and
 650 $\{m+1, \dots, n\}$ (block 2), where $m \ll n$, and focus on the first m columns of \mathbf{Z} for the update. We partition
 651 \mathbf{Z} , \mathbf{C} , \mathbf{S} , and \mathbf{K}^{-1} as follows:

$$\mathbf{Z} = [\mathbf{Z}_1 \quad \mathbf{Z}_2], \quad \mathbf{K}^{-1} = \begin{bmatrix} \mathbf{K}_{11}^{-1} & \mathbf{K}_{12}^{-1} \\ \mathbf{K}_{12}^{-\top} & \mathbf{K}_{22}^{-1} \end{bmatrix}, \quad \boldsymbol{\gamma} = \begin{bmatrix} c(\mathbf{z}_1, \mathbf{z}_1) & \cdots & c(\mathbf{z}_1, \mathbf{z}_m) \\ \vdots & \ddots & \vdots \\ c(\mathbf{z}_m, \mathbf{z}_1) & \cdots & c(\mathbf{z}_m, \mathbf{z}_m) \end{bmatrix}, \quad (27)$$

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{12}^\top & \mathbf{S}_{22} \end{bmatrix}, \quad \mathbf{C} = \begin{bmatrix} \boldsymbol{\gamma} & \boldsymbol{\beta} \\ \boldsymbol{\beta}^\top & \mathbf{C}_{22} \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} c(\mathbf{z}_1, \mathbf{z}_{m+1}) & \cdots & c(\mathbf{z}_1, \mathbf{z}_n) \\ \vdots & \ddots & \vdots \\ c(\mathbf{z}_m, \mathbf{z}_{m+1}) & \cdots & c(\mathbf{z}_m, \mathbf{z}_n) \end{bmatrix}.$$

652 Here, subscripts represent the block indices, so \mathbf{Z}_1 and \mathbf{Z}_2 are the first m columns and the last $n - m$
 653 columns of \mathbf{Z} and subscript 12 indicates the block matrix across the first m variables and the last $n - m$
 654 variables. Only $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ contain the active variables in \mathbf{Z}_1 to optimize. \mathbf{Z}_2 is fixed.

655 Applying the block representation to eq. 15, we get the block PLS objective function for solving \mathbf{Z}_1 :

$$\begin{aligned} l_{\text{PLS}}(\mathbf{Z}_1) &= \text{tr} [(\mathbf{S}_{11} - \boldsymbol{\gamma} - \sigma^2 \mathbf{I}_n)(\mathbf{S}_{11} - \boldsymbol{\gamma} - \sigma^2 \mathbf{I}_n)^T + 2(\mathbf{S}_{12} - \boldsymbol{\beta})(\mathbf{S}_{12} - \boldsymbol{\beta})^T] \\ &\quad + \text{tr} [(\mathbf{Z}_1 - \mathbf{B}\mathbf{X}_1)\mathbf{K}_{11}^{-1}(\mathbf{Z}_1 - \mathbf{B}\mathbf{X}_1)^T + 2(\mathbf{Z}_1 - \mathbf{B}\mathbf{X}_1)\mathbf{K}_{12}^{-1}(\mathbf{Z}_2 - \mathbf{B}\mathbf{X}_2)^T], \end{aligned} \quad (28)$$

656 where \mathbf{X}_1 and \mathbf{X}_2 are the first m columns and the last $n - m$ columns of \mathbf{X} .

657 To formulate the block MAP estimator for \mathbf{Z}_1 , we first apply the block matrix inversion to \mathbf{C} ,

$$\mathbf{C}^{-1} = \begin{bmatrix} (\boldsymbol{\gamma} - \boldsymbol{\beta}\mathbf{C}_{22}^{-1}\boldsymbol{\beta}^T)^{-1}, & -\boldsymbol{\gamma}^{-1}\boldsymbol{\beta}(\mathbf{C}_{22} - \boldsymbol{\beta}^T\boldsymbol{\gamma}^{-1}\boldsymbol{\beta})^{-1} \\ -\mathbf{C}_{22}^{-1}\boldsymbol{\beta}^T(\boldsymbol{\gamma} - \boldsymbol{\beta}\mathbf{C}_{22}^{-1}\boldsymbol{\beta}^T)^{-1}, & (\mathbf{C}_{22} - \boldsymbol{\beta}^T\boldsymbol{\gamma}^{-1}\boldsymbol{\beta})^{-1} \end{bmatrix}. \quad (29)$$

658 Incorporating this matrix inversion into the MAP estimator, we get the objective function

$$\begin{aligned} l_{\text{MAP}}(\mathbf{Z}_1) &= \log|\boldsymbol{\gamma} + \sigma^2 \mathbf{I}_n| + \text{tr} [\boldsymbol{\beta}\mathbf{C}_{22}^{-1}\mathbf{S}_{22}\mathbf{C}_{22}^{-1}\boldsymbol{\beta}^T(\boldsymbol{\gamma} + \sigma^2 \mathbf{I}_n)^{-1}] \\ &\quad - 2\text{tr} [\mathbf{S}_{12}\mathbf{C}_{22}^{-1}\boldsymbol{\beta}^T(\boldsymbol{\gamma} + \sigma^2 \mathbf{I}_n)^{-1}] + \text{tr} [\mathbf{S}_{11}(\boldsymbol{\gamma} + \sigma^2 \mathbf{I}_n)^{-1}] \\ &\quad + \text{tr} [(\mathbf{Z}_1 - \mathbf{B}\mathbf{X}_1)\mathbf{K}_{11}^{-1}(\mathbf{Z}_1 - \mathbf{B}\mathbf{X}_1)^T + 2(\mathbf{Z}_1 - \mathbf{B}\mathbf{X}_1)\mathbf{K}_{12}^{-1}(\mathbf{Z}_2 - \mathbf{B}\mathbf{X}_2)^T]. \end{aligned} \quad (30)$$

659 The time complexity of the block PLS estimator is $O(nm^2)$ per iteration, where m is the size of the block,
660 and the time complexity of the MAP estimator is $O(n^2m)$ per iteration. For comparison, greedy gradient
661 descent has complexity $O(n^3)$ per iteration. In the experiments, the block MAP estimator for \mathbf{Z} with the
662 block PLS estimator as a warm start is a practical optimization approach.

663 We now describe the block coordinate descent (BCD) update. We assume that the voxel indices $\{1, \dots, n\}$
664 are divided into k blocks $\{I_j\}_{j=1}^k$, where I_j is the set of indices corresponding to the columns of \mathbf{Z} in the
665 j 'th block. Denote I_j columns of \mathbf{Z} to be \mathbf{Z}_{I_j} . We cluster indices into blocks based on the spatial locations
666 of the voxels and assume smooth measurements for nearby voxels. Thus, the size of a block should be at
667 least one length scale of the region defined by the kernel in \mathbf{x} space to encourage dependencies among
668 neighboring voxels. This smoothness assumption leads to a block-wise but not an element-wise update,
669 which separates our BCD method from Informative Vector Machine (IVM) [48]. At each iteration t , we
670 choose a nonempty index subset $I \in \{I_j\}_{j=1}^k$. Then the objective functions $l(\mathbf{Z}_I)$ in eq. 28 and 30 are
671 optimized w.r.t. \mathbf{Z}_I via L-BFGS [49]. After the t 'th iteration, \mathbf{Z}_I is updated, holding all other blocks fixed.

672 For high-dimensional non-convex problems, a good initialization is essential to finding a good optimum.
673 Here we describe an effective and tractable initialization for a single block of \mathbf{Z} , inspired by Laplacian
674 eigenmaps.

675 The Laplacian eigenmap (LE) algorithm is a popular dimensionality reduction method that solves a
676 generalized eigendecomposition [50]. LE defines a neighborhood graph on the data $\{\mathbf{y}_i \in \mathbb{R}^T\}_{i=1}^n$, such
677 as k nearest neighbors or an ϵ -ball graph, and weighs each graph edge $\mathbf{y}_i \sim \mathbf{y}_j$ by a symmetric affinity
678 function $V(\mathbf{y}_i, \mathbf{y}_j) = v_{ij}$, typically Gaussian: $v_{ij} = \exp\left(-\frac{\|\mathbf{y}_i - \mathbf{y}_j\|^2}{2s^2}\right)$ with s the length scale. Given this
679 weighted graph, LE seeks latent points $\{\mathbf{z}_i \in \mathbb{R}^d\}_{i=1}^n$ that are solutions to the optimization problem

$$\min_{\mathbf{Z}} \text{tr}(\mathbf{Z}^T \mathbf{L} \mathbf{Z}) \quad \text{s.t.} \quad \mathbf{Z}^T \mathbf{D} \mathbf{Z} = \mathbf{I}, \quad \mathbf{Z}^T \mathbf{D} \mathbf{1} = \mathbf{I}, \quad (31)$$

680 with the symmetric affinity matrix $\mathbf{V} \in \mathbb{R}^{n \times n}$, the degree matrix $\mathbf{D} = \text{diag}(\sum_{i=1}^n v_{ii}) \in \mathbb{R}^{n \times n}$, the graph
681 Laplacian matrix $\mathbf{L} = \mathbf{D} - \mathbf{V} \in \mathbb{R}^{n \times n}$, and $\mathbf{1} = [1, \dots, 1]^T$. Constraints eliminate the two trivial solutions
682 $\mathbf{Z} = \mathbf{0}$ by setting an arbitrary scale and $\mathbf{z}_1 = \dots = \mathbf{z}_n$ by removing $\mathbf{1}$, which is an eigenvector of \mathbf{L}
683 associated with a zero eigenvalue.

684 Following the previous two-block example, let \mathbf{Z} be partitioned into \mathbf{Z}_1 and \mathbf{Z}_2 . To update \mathbf{Z}_1 given \mathbf{Z}_2 , the
685 objective function is:

$$\min_{\mathbf{Z}_1} \text{tr} \left([\mathbf{Z}_1 \quad \mathbf{Z}_2] \begin{bmatrix} \mathbf{L}_{11} & \mathbf{L}_{12} \\ \mathbf{L}_{12}^T & \mathbf{L}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{Z}_1^T \\ \mathbf{Z}_2^T \end{bmatrix} \right). \quad (32)$$

686 We don't need to use the constraints from eq. 31 because the trivial solutions are removed given \mathbf{Z}_2 . The
687 solution is

$$\tilde{\mathbf{Z}}_1 = -\mathbf{Z}_2 \mathbf{L}_{12}^\top \mathbf{L}_{11}^{-1}. \quad (33)$$

688 Given the current latent embeddings for all other coordinates, the algorithm seeks the best latent embed-
689 ding of the unknown dimensions. Because of the computational efficiency of this approach, we use $\tilde{\mathbf{Z}}_1$ to
690 initialize the BCD algorithm for each iteration and \mathbf{Z}_2 collapses all the fixed blocks.

691 In addition, if \mathbf{V} has a Gaussian affinity function, the latent embedding \mathbf{Z} is mapped from the observation \mathbf{Y}
692 with a radial basis function nonlinearity. For covariance estimation, we enforce the resemblance between
693 another radial basis function (RBF) kernel on \mathbf{Z} and the sample covariance of \mathbf{Y} . We are essentially trying
694 to map the observation space to itself with double layers of exponential transformations, which would result
695 in a bad latent embedding for initialization. Therefore, instead of using a Gaussian function for the weights
696 \mathbf{V} , we use a function that inverts the RBF nonlinearity on a covariance matrix, defined as $v_{ij} = -f(\mathbf{y}_i \mathbf{y}_j^\top)$.
697 Here, $f(x) = \text{sign}(x) \log(|x| + 1)$ is the log-modulus transformation [51], which distributes the magnitude
698 of the data while preserving the sign of the data in order to control against negative covariance values
699 when taking the logarithm.

700 Hyperparameter estimation for the brain kernel

701 Our model includes five hyperparameters: $\{\delta, r, \mathbf{B}, \rho, \sigma^2\}$. We set these parameters as follows.

702 **Estimate δ :** δ is the length scale of the GP kernel mapping from coordinate space to latent space. We can
703 optimize this parameter by taking the derivative of the GP log-likelihood w.r.t. δ ; this involves inverting the
704 kernel matrix with an $O(n^3)$ cost, which is computationally infeasible here. An inducing-point method [18]
705 introduced extra inducing variables to optimize, which further increased the computational burden. Instead,
706 we use a scalable spectral formulation for learning the length scale δ of the GP kernel.

707 For a stationary Gaussian process $f \sim \mathcal{GP}(m, k)$, when f has a high degree of smoothness, the prior
708 covariance \mathbf{K} becomes approximately low rank, meaning that it has a small number of non-negligible
709 eigenvalues. Because the kernel function for \mathbf{x} space is shift invariant, the eigenspectrum of \mathbf{K} has a
710 diagonal representation in the Fourier domain, a consequence of Bochner's theorem [52–54],

$$g \sim \mathcal{N}(\alpha(\omega), \Sigma(\omega)), \quad (34)$$

711 where g is the Fourier transform of f , ω is the frequency, $\alpha(\omega)$ is the Fourier transform of the mean function
712 $m(\mathbf{x})$, and $\Sigma = \text{diag}(s(\omega))$ is diagonal. This means that Fourier components are a priori independent,
713 with prior variance

$$s(\omega) = (2\pi\delta^2)^{h/2} r \exp(-2\pi^2\delta^2\omega^2), \quad (35)$$

714 where h is the dimension of the input \mathbf{x} . Without loss of generality, we assume the size of the spectral
715 domain for each input dimension is w , and thus $\alpha = \alpha(\omega) \in \mathbb{R}^{w^h}$ and $\Sigma = \Sigma(\omega) \in \mathbb{R}^{w^h \times w^h}$ given a
716 spectral representation ω . The mapping between $f(\mathbf{x})$ and its Fourier transform $g(\omega)$ is then

$$f(\mathbf{x}) = \sum_j e^{2\pi i \mathbf{x}^\top \omega_j} g(\omega_j) = \mathbf{e}^\top g(\omega), \quad (36)$$

717 where \mathbf{e} is a column vector with entries $e^{2\pi i \mathbf{x}^\top \omega_j}$ on the j 'th position for the j 'th frequency. Let $\mathbf{e}_i \in \mathbb{R}^{w^h}$
718 denote the Fourier vector for \mathbf{x}_i , then we can further define $\mathcal{B} = (\mathbf{e}_1, \dots, \mathbf{e}_n) \in \mathbb{R}^{w^h \times n}$ to be the
719 Fourier matrix for an input coordinate matrix \mathbf{X} . Let $\mathbf{z}_j \in \mathbb{R}^n$ denote the j 'th latent embedding of \mathbf{X}
720 and $j \in \{1, \dots, d\}$. Note that \mathbf{z}_j is the j 'th row of the matrix \mathbf{Z} . We can write $\mathbf{z}_j = \mathbf{g}_j^\top \mathcal{B}$, where

721 $\mathbf{g}_j \in \mathbb{R}^{w^h} \sim \mathcal{N}(\boldsymbol{\alpha}_j, \Sigma)$. This implies that each latent embedding deterministically depends on a unique
 722 spectral function. These spectral functions are all sampled from multivariate Gaussian priors with different
 723 mean functions but the same covariance function in the spectral domain. Bringing the Fourier matrix
 724 \mathcal{B} into the Gaussian prior (eq. 34), we derive the Gaussian prior over the latent \mathbf{z}_j expressed with the
 725 spectral formulation as

$$\mathbf{z}_j \sim \mathcal{N}(\boldsymbol{\alpha}_j^\top \mathcal{B}, \mathcal{B}^\top \Sigma \mathcal{B}). \quad (37)$$

726 We can use this representation to optimize δ . If δ is in a large-value regime, we can ignore Fourier
 727 components above a certain high-frequency cutoff which leads to a lower-dimensional ω and a lower-
 728 dimensional optimization problem. Because $h = 3$ in fMRI analyses, we can control w to be small enough
 729 so that w^h is tractable relative to the large n , and \mathcal{B} is a manageable high-dimensional matrix. Because
 730 we generally do not assume uniform gridding of coordinates \mathbf{X} , \mathcal{B} is a non-uniform DFT (not orthogonal).
 731 Although the Kronecker tricks cannot be used for scaling up with non-uniform DFT, we can employ block
 732 matrix inverse lemma to transform the spectral formulation into a lower-dimensional problem.

733 Then the standard negative GP log-likelihood to optimize for δ is

$$l(\delta) = \text{tr}[(\mathbf{Z} - \mathbf{B}\mathbf{X})(\mathbf{K} + \epsilon\mathbf{I})^{-1}(\mathbf{Z} - \mathbf{B}\mathbf{X})^\top] + \log|\mathbf{K} + \epsilon\mathbf{I}|, \quad (38)$$

734 where ϵ is a small noise variance in the latent space to avoid ill-conditions for \mathbf{K} . Let $\mathbf{A} = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_d)^\top \in$
 735 $\mathbb{R}^{d \times w^h}$. With the spectral representation, eq. 38 can be re-written as

$$l(\delta) = \text{tr}[(\mathbf{Z} - \mathbf{A}\mathcal{B})(\mathcal{B}^\top \Sigma \mathcal{B} + \epsilon\mathbf{I})^{-1}(\mathbf{Z} - \mathbf{A}\mathcal{B})^\top] + \log|\mathcal{B}^\top \Sigma \mathcal{B} + \epsilon\mathbf{I}|. \quad (39)$$

736 Calculating $(\mathcal{B}^\top \Sigma \mathcal{B} + \epsilon\mathbf{I})^{-1}$ is still computationally expensive, but with the spectral factorization, we are
 737 able to use the block matrix inverse lemma as in eq. 29,

$$\begin{aligned} & (\mathcal{B}^\top \Sigma \mathcal{B} + \epsilon\mathbf{I})^{-1} = \\ & \left[\mathbf{I} - \mathcal{B}_1^\top (\mathbf{P}_2^{-1} + \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \mathcal{B}_1, -\mathcal{B}_1^\top (\Sigma \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \Sigma (\mathcal{B}_{22} \Sigma - \mathcal{B}_{22} \Sigma (\mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_2 \right. \\ & \left. -\mathcal{B}_2^\top (\Sigma \mathcal{B}_{22} + \epsilon\mathbf{I})^{-1} \Sigma (\mathcal{B}_{11} \Sigma - \mathcal{B}_{11} \Sigma (\mathcal{B}_{22} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_{22} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_1, \mathbf{I} - \mathcal{B}_2^\top (\mathbf{P}_1^{-1} + \mathcal{B}_{22} + \epsilon\mathbf{I})^{-1} \mathcal{B}_2 \right] \end{aligned} \quad (40)$$

738 where \mathcal{B}_1 and \mathcal{B}_2 correspond to the Fourier bases for \mathbf{X}_1 and \mathbf{X}_2 respectively. Note that neither \mathcal{B}_1
 739 nor \mathcal{B}_2 is invertible. Moreover, $\mathcal{B}_{11} = \mathcal{B}_1 \mathcal{B}_1^\top$, $\mathcal{B}_{22} = \mathcal{B}_2 \mathcal{B}_2^\top$, $\mathbf{P}_1 = \Sigma - \Sigma (\mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_{11} \Sigma$ and
 740 $\mathbf{P}_2 = \Sigma - \Sigma (\mathcal{B}_{22} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_{22} \Sigma$. All matrices have size $w^h \times w^h$, which is tractable to invert. We know
 741 that the spectral expression of \mathbf{K}^{-1} is $(\mathcal{B}^\top \Sigma \mathcal{B} + \epsilon\mathbf{I})^{-1}$. We can also express \mathbf{K}_{11}^{-1} and \mathbf{K}_{12}^{-1} with the
 742 spectral formulation as

$$\begin{aligned} \mathbf{K}_{11}^{-1} &= \mathbf{I} - \mathcal{B}_1^\top (\mathbf{P}_2^{-1} + \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \mathcal{B}_1, \\ \mathbf{K}_{12}^{-1} &= -\mathcal{B}_1^\top (\Sigma \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \Sigma (\mathcal{B}_{22} \Sigma - \mathcal{B}_{22} \Sigma (\mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_2. \end{aligned} \quad (41)$$

743 Consequently, the negative GP log-likelihood w.r.t. δ in the block form is represented as

$$\begin{aligned} l(\delta) &= (\mathbf{Z}_1 - \mathbf{A}\mathcal{B}_1)(\mathbf{I} - \mathcal{B}_1^\top (\mathbf{P}_2^{-1} + \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \mathcal{B}_1)(\mathbf{Z}_1 - \mathbf{A}\mathcal{B}_1)^\top \\ &+ (\mathbf{Z}_2 - \mathbf{A}\mathcal{B}_2)(\mathbf{I} - \mathcal{B}_2^\top (\mathbf{P}_1^{-1} + \mathcal{B}_{22} + \epsilon\mathbf{I})^{-1} \mathcal{B}_2)(\mathbf{Z}_2 - \mathbf{A}\mathcal{B}_2)^\top \\ &+ \log|\Sigma \mathcal{B}_{22} - \Sigma \mathcal{B}_{11} (\Sigma \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \Sigma \mathcal{B}_{22} + \epsilon\mathbf{I}| + \log|\Sigma \mathcal{B}_{11} + \epsilon\mathbf{I}| \\ &- 2(\mathbf{Z}_1 - \mathbf{A}\mathcal{B}_1)\mathcal{B}_1^\top (\Sigma \mathcal{B}_{11} + \epsilon\mathbf{I})^{-1} \Sigma (\mathcal{B}_{22} \Sigma - \mathcal{B}_{22} \Sigma (\mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \\ &\cdot \mathcal{B}_{11} \Sigma + \epsilon\mathbf{I})^{-1} \mathcal{B}_2(\mathbf{Z}_2 - \mathbf{A}\mathcal{B}_2)^\top. \end{aligned} \quad (42)$$

744 With eq. 42, the computational cost is reduced to $\max(O(nw^h), O(w^{3h}))$, where w is the dimension of
 745 the spectral form per input dimension and h is the number of input dimensions. This complexity is much
 746 smaller than $O(n^3)$ when $n \gg w$ and $h \leq 3$. Another benefit of this formulation is that δ only exists in the
 747 diagonal of Σ (eq. 35), which makes the estimation straightforward.

Algorithm 1 Block Coordinate Descent Algorithm for Brain Kernel

Input: sample covariance \mathbf{S} , voxel coordinate \mathbf{X}

Output: latent variable \mathbf{Z} , linear projection \mathbf{B} , length scale δ , marginal variance r

Generate the index set $\{I_j\}_{j=1}^k$ given \mathbf{X}

Randomly initialize \mathbf{Z}

for $t = 1, 2, \dots$ **do**

Pick $I \in \{I_j\}_{j=1}^k$, and solve eq. 33 to get an initialization for $\mathbf{Z}_I^{(t)}$

Find the optimal $\mathbf{Z}_I^{(t)}$ and $\mathbf{B}^{(t)}$ by solving eq. 28 or 30

Estimate the optimal length scale δ^* and the optimal marginal variance r^* by solving eq. 42

Update $\mathbf{K}_{11}^{-1(t)}$ and $\mathbf{K}_{12}^{-1(t)}$ given δ^* and r^* according to eq. 41

end for

748 **Estimate r :** r determines the scale of the latent embedding \mathbf{Z} . Since r also exists on the diagonal of Σ
 749 (eq. 35), the optimization for r can be combined with learning δ using the same spectral representation in
 750 eq. 42.

751 **Estimate ρ :** ρ is the marginal variance of the covariance function. We assume that the data has been
 752 normalized to have zero mean and variance one. Thus we set $\rho = 1$.

753 **Estimate \mathbf{B} :** \mathbf{B} is the linear projection of the mean function. We can estimate \mathbf{B} jointly with \mathbf{Z} in each
 754 BCD iteration by optimizing the same objective function (eq. 28 and 30).

755 **Estimate σ^2 :** σ^2 is the observation noise variance. We estimate σ^2 using the eigenvalues of the sample
 756 covariance \mathbf{S} [55].

757

758 Algorithm 1 describes the complete BCD algorithm for the brain kernel. The convergence of the BCD
 759 algorithm (without parameter estimation) to a stationary point is addressed in the theoretical results in
 760 previous work [22]. There, a general block-coordinate-descent approach is analyzed to solve minimization
 761 problems of the form $F(x) = f(x) + \lambda h(x)$, which is composed of the sum of a smooth function $f(\cdot)$ and
 762 a separable convex function $h(\cdot)$, in our case $h(x) = 0$. According to Part (e) of Theorem 4.1 in [22], if I_t
 763 at the t 'th iteration is chosen by the generalized Gauss-Seidel rule,

$$\bigcup_{i=0,1,\dots,T-1} J_{i+t} \supseteq \mathcal{V} \quad \forall t = 1, 2, \dots, T, \quad (43)$$

764 which is necessary to ensure that all variables are updated every T steps and \mathcal{V} is the set of all variables,
 765 then each coordinate-wise minimum point of $\{\mathbf{Z}_{I_t}\}$ is a stationary point of $f(\mathbf{Z})$.

766 **Hyperparameter optimization and latent variable inference for Bayesian factor
 767 analysis**

Our goal was to infer the latent variables \mathbf{L} and \mathbf{F} and the hyperparameters for \mathbf{C} and σ^2 , denoted as θ .
 We assumed that

$$\mathbf{X} = \mathbf{LF} + \epsilon \sim \mathcal{N}(\mathbf{LF}, \sigma^2 \mathbf{I}), \quad (44)$$

$$\mathbf{L}_j \sim \mathcal{N}(\mathbf{0}, \mathbf{C}), \quad \text{for the } j\text{th column}, \quad (45)$$

and

$$\mathbf{F}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad \text{for the } i\text{th column}. \quad (46)$$

We first aimed at estimating the hyperparameters θ by marginalizing over \mathbf{L} and \mathbf{F} . The marginal distribution for \mathbf{X} is

$$p(\mathbf{X}) = \int \mathcal{N}(\mathbf{LF}, \sigma^2 \mathbf{I}) \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) \prod_{i=1}^T \mathcal{N}(\mathbf{F}_i | \mathbf{0}, \mathbf{I}) d\mathbf{L}_1 \dots d\mathbf{L}_K d\mathbf{F}_1 \dots d\mathbf{F}_T, \quad (47)$$

which is intractable. However, we could match the second order moment of \mathbf{X} to the sample covariance $\mathbf{S} = \tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$ where $\tilde{\mathbf{X}}$ is the standardized data sample. The second order moment is derived as follows

$$\begin{aligned} \mathbb{E}[\mathbf{XX}^\top] &= \mathbb{E}[\mathbf{LFF}^\top \mathbf{L}^\top] + \sigma^2 T \mathbf{I} \\ &= \int \int \mathbf{LFF}^\top \mathbf{L}^\top \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) \prod_{i=1}^T \mathcal{N}(\mathbf{F}_i | \mathbf{0}, \mathbf{I}) d\mathbf{L}_1 \dots d\mathbf{L}_K d\mathbf{F}_1 \dots d\mathbf{F}_T + \sigma^2 T \mathbf{I} \\ &= \int \mathbf{L} \left[\int \mathbf{FF}^\top \prod_{i=1}^T \mathcal{N}(\mathbf{F}_i | \mathbf{0}, \mathbf{I}) d\mathbf{F}_1 \dots d\mathbf{F}_T \right] \mathbf{L}^\top \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) d\mathbf{L}_1 \dots d\mathbf{L}_K + \sigma^2 T \mathbf{I} \\ &= \int \mathbf{L} \left[\int \sum_{i=1}^T \mathbf{F}_i \mathbf{F}_i^\top \prod_{i=1}^T \mathcal{N}(\mathbf{F}_i | \mathbf{0}, \mathbf{I}) d\mathbf{F}_1 \dots d\mathbf{F}_T \right] \mathbf{L}^\top \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) d\mathbf{L}_1 \dots d\mathbf{L}_K + \sigma^2 T \mathbf{I} \\ &= \int \mathbf{L} \left[\sum_{i=1}^T \int \mathbf{F}_i \mathbf{F}_i^\top \mathcal{N}(\mathbf{F}_i | \mathbf{0}, \mathbf{I}) d\mathbf{F}_i \right] \mathbf{L}^\top \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) d\mathbf{L}_1 \dots d\mathbf{L}_K + \sigma^2 T \mathbf{I} \\ &= T \int \mathbf{LL}^\top \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) d\mathbf{L}_1 \dots d\mathbf{L}_K + \sigma^2 T \mathbf{I} \\ &= T \sum_{j=1}^K \int \mathbf{L}_j \mathbf{L}_j^\top \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}) d\mathbf{L}_j + \sigma^2 T \mathbf{I} \\ &= TK\mathbf{C} + \sigma^2 T \mathbf{I}. \end{aligned} \quad (48)$$

⁷⁶⁸ By matching eq. 48 to the sample covariance \mathbf{S} , we were able to estimate the hyperparameters in \mathbf{C} and
⁷⁶⁹ σ^2 .

Next, we fixed the estimated hyperparameters and inferred \mathbf{L} . Marginalizing over \mathbf{F} we got

$$p(\mathbf{X}|\mathbf{L}) = \mathcal{N}(\mathbf{0}, \mathbf{LL}^\top + \sigma^2 \mathbf{I}). \quad (49)$$

Then, we obtained the joint distribution between \mathbf{X} and \mathbf{L} as

$$p(\mathbf{X}, \mathbf{L}) = \mathcal{N}(\mathbf{0}, \mathbf{LL}^\top + \sigma^2 \mathbf{I}) \prod_{j=1}^K \mathcal{N}(\mathbf{L}_j | \mathbf{0}, \mathbf{C}), \quad (50)$$

whose log likelihood is written as

$$\mathcal{L}(\mathbf{L}) = -\frac{1}{2} \log |\mathbf{LL}^\top + \sigma^2 \mathbf{I}| - \frac{1}{2T} \text{Tr} (\mathbf{X}^\top (\mathbf{LL}^\top + \sigma^2 \mathbf{I})^{-1} \mathbf{X}) - \frac{1}{2K} \text{Tr} (\mathbf{L}^\top \mathbf{C}^{-1} \mathbf{L}). \quad (51)$$

We maximized eq. 51 w.r.t. \mathbf{L} . To speed up the optimization process, we initialized \mathbf{L} via finding a closed-form solution approximately maximizing eq. 51. Denoting $\mathbf{LL}^\top + \sigma^2 \mathbf{I}$ as \mathbf{Q} , we rewrote the log likelihood as

$$\mathcal{L}(\mathbf{Q}) = -\frac{1}{2} \log |\mathbf{Q}| - \frac{1}{2T} \text{Tr} (\mathbf{Q}^{-1} \mathbf{XX}^\top) - \frac{1}{2K} \text{Tr} (\mathbf{C}^{-1} \mathbf{Q}), \quad (52)$$

whose derivative is

$$\frac{\partial \mathcal{L}(\mathbf{Q})}{\partial \mathbf{Q}} = -\frac{1}{2}\mathbf{Q}^{-1} + \frac{1}{2T}\mathbf{Q}^{-1}\mathbf{X}\mathbf{X}^\top\mathbf{Q}^{-1} - \frac{1}{2K}\mathbf{C}^{-1}. \quad (53)$$

Setting the derivative to be 0, we arrived at

$$\begin{aligned} \frac{1}{K}\mathbf{Q}\mathbf{C}^{-1}\mathbf{Q} + \mathbf{Q} &= \frac{1}{T}\mathbf{X}\mathbf{X}^\top \\ \Downarrow \\ (\mathbf{Q} + \frac{1}{2}K\mathbf{C})\frac{1}{K}\mathbf{C}^{-1}(\mathbf{Q} + \frac{1}{2}K\mathbf{C}) &= \frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C}. \end{aligned} \quad (54)$$

⁷⁷⁰ We needed to decompose the right-hand side into the multiplication of a symmetric matrix, \mathbf{C}^{-1} and the
⁷⁷¹ same symmetric matrix. Here is our solution:

- ⁷⁷² • Let \mathbf{P} denote the Cholesky decomposition of $\frac{1}{K}\mathbf{C}^{-1}$, i.e., $\frac{1}{K}\mathbf{C}^{-1} = \mathbf{P}\mathbf{P}^\top$.
- ⁷⁷³ • Denote $\mathbf{Q} + \frac{1}{2}K\mathbf{C} = \mathbf{B}\mathbf{B}^\top$ and let $\mathbf{B} = \mathbf{P}^{-\top}\mathbf{A}$ where \mathbf{A} is an unknown square matrix.
- Then we can rewrite eq. 54 as

$$\begin{aligned} \frac{1}{K}\mathbf{B}\mathbf{B}^\top\mathbf{C}^{-1}\mathbf{B}\mathbf{B}^\top &= \frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C} \\ \Rightarrow \mathbf{P}^{-\top}\mathbf{A}\mathbf{A}^\top\mathbf{P}^{-1}\mathbf{P}\mathbf{P}^\top\mathbf{P}^{-\top}\mathbf{A}\mathbf{A}^\top\mathbf{P}^{-1} &= \frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C} \\ \Rightarrow \mathbf{P}^{-\top}\mathbf{A}\mathbf{A}^\top\mathbf{A}\mathbf{A}^\top\mathbf{P}^{-1} &= \frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C} \\ \Rightarrow \mathbf{A}\mathbf{A}^\top\mathbf{A}\mathbf{A}^\top &= \mathbf{P}^\top(\frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C})\mathbf{P}. \end{aligned}$$

- ⁷⁷⁴ • Next, we represent \mathbf{A} with its singular value decomposition (SVD), i.e., $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^\top$.
- Finally we can rewrite eq. 54 as

$$\begin{aligned} \mathbf{A}\mathbf{A}^\top\mathbf{A}\mathbf{A}^\top &= \mathbf{P}^\top(\frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C})\mathbf{P} \\ \Rightarrow \mathbf{U}\mathbf{S}\mathbf{V}^\top\mathbf{V}\mathbf{S}\mathbf{U}^\top\mathbf{U}\mathbf{S}\mathbf{V}^\top\mathbf{V}\mathbf{S}\mathbf{U}^\top &= \mathbf{P}^\top(\frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C})\mathbf{P} \\ \Rightarrow \mathbf{U}\mathbf{S}^4\mathbf{U}^\top &= \mathbf{P}^\top(\frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C})\mathbf{P}. \end{aligned} \quad (55)$$

- ⁷⁷⁵ • Therefore, to solve eq. 54, we first factorize $\mathbf{P}^\top(\frac{1}{T}\mathbf{X}\mathbf{X}^\top + \frac{1}{4}K\mathbf{C})\mathbf{P}$ using SVD to obtain \mathbf{U} and \mathbf{S} according to eq. 55.
- ⁷⁷⁶ • Then $\mathbf{Q} = \mathbf{B}\mathbf{B}^\top - \frac{1}{2}K\mathbf{C} = \mathbf{P}^{-\top}\mathbf{A}\mathbf{A}^\top\mathbf{P}^{-1} - \frac{1}{2}K\mathbf{C} = \mathbf{P}^{-\top}\mathbf{U}\mathbf{S}^2\mathbf{U}^\top\mathbf{P}^{-1} - \frac{1}{2}K\mathbf{C}$.
- ⁷⁷⁷ • Ultimately, \mathbf{L} can be obtained via $\mathbf{L} = \mathbf{W}_{:,K}\mathbf{D}_{:K,:K}^{\frac{1}{2}}$ where $\mathbf{W}\mathbf{D}\mathbf{W}^\top$ is the eigen-decomposition of $\mathbf{Q} - \sigma^2\mathbf{I}$.

Given the above procedure, we were able to find an ideal initialization for \mathbf{L} which made the optimization of eq. 51 much easier. Conditioned on the optimal \mathbf{L}^* , we turned to inferring the optimal \mathbf{F}^* via maximizing

$$\mathcal{L}(\mathbf{F}) = \log \mathcal{N}(\mathbf{X}|\mathbf{L}^*\mathbf{F}, \sigma^2\mathbf{I}) \prod_{i=1}^T \mathcal{N}(\mathbf{F}_i|\mathbf{0}, \mathbf{I}), \quad (56)$$

780 which has a closed-form solution, i.e., $\mathbf{F}^* = (\mathbf{L}^{*\top} \mathbf{L}^* + \sigma^2 \mathbf{I})^{-1} \mathbf{L}^{*\top} \mathbf{X}$.

781 Note that in order to obtain $\mathbf{L} = \mathbf{W}_{:,K} \mathbf{D}_{:K,:K}^{\frac{1}{2}}$, we need to guarantee that we are able to find a positive
782 semi-definite (p.s.d.) matrix $\mathbf{Q} - \sigma^2 \mathbf{I}$. Here are the lemma and the theorem:

783 **Lemma.** *If \mathbf{A} is p.s.d., and \mathbf{B} is symmetric p.s.d., then \mathbf{AB} is also p.s.d..*

784 *Proof.* If \mathbf{A} and \mathbf{B} are both p.s.d. and \mathbf{B} is also symmetric, then suppose λ is an eigenvalue of \mathbf{AB} with
785 corresponding eigenvector $\mathbf{x} \neq 0$, i.e., $\mathbf{ABx} = \lambda \mathbf{x}$. Then $\mathbf{BABx} = \lambda \mathbf{Bx}$ and so $\mathbf{x}^\top \mathbf{BABx} = \lambda \mathbf{x}^\top \mathbf{Bx}$. Both
786 It is not hard to check that \mathbf{BAB} will also be p.s.d.. For \mathbf{x} s.t. $\mathbf{x}^\top \mathbf{Bx} \neq 0$, we have $\lambda = \frac{\mathbf{x}^\top \mathbf{BABx}}{\mathbf{x}^\top \mathbf{Bx}}$. Both
787 the numerator and the denominator are non-negative values, therefore $\lambda \geq 0$. For \mathbf{x} s.t. $\mathbf{x}^\top \mathbf{Bx} = 0$,
788 we assume $\mathbf{x} = \mathbf{Ve}$ where $\mathbf{B} = \mathbf{VDV}^\top$ is the eigen-decomposition and \mathbf{e} is the linear weight vector.
789 Then $\mathbf{x}^\top \mathbf{Bx} = \mathbf{e}^\top \mathbf{V}^\top \mathbf{BVe} = \mathbf{e}^\top \mathbf{V}^\top \mathbf{VDe} = \mathbf{e}^\top \mathbf{De} = 0$. Since \mathbf{D} is a diagonal matrix with non-
790 negative values, \mathbf{e} should have zero elements corresponding to the non-zero eigenvalues. Therefore
791 \mathbf{x} is a linear combination of eigenvectors of \mathbf{B} whose eigenvalues are zero, i.e., $\mathbf{x} = \mathbf{V}_0 \mathbf{e}$ where \mathbf{V}_0
792 contains all zero eigenvectors with $\mathbf{D}_0 = \mathbf{0}$ and \mathbf{e} has no zero elements. Following that, we have
793 $\mathbf{Bx} = \mathbf{BV}_0 \mathbf{e} = \mathbf{V}_0 \mathbf{D}_0 \mathbf{e} = \mathbf{0} \Rightarrow \mathbf{ABx} = \lambda_B \mathbf{Ax} = \mathbf{0} \Rightarrow \lambda = 0$. Based on the derivation, we arrive at the
794 conclusion that \mathbf{AB} has non-negative eigenvalues thus \mathbf{AB} is p.s.d.. \square

795 **Theorem.** *If $K\mathbf{C} - \sigma^2 \mathbf{I}$ is p.s.d., then $\mathbf{Q} - \sigma^2 \mathbf{I}$ is p.s.d. based on the Lemma.*

Proof.

$$\begin{aligned}
 & K\mathbf{C} \succeq \sigma^2 \mathbf{I} \\
 & \Rightarrow K\mathbf{C} \frac{1}{T} \mathbf{XX}^\top \succeq \sigma^2 \frac{1}{T} \mathbf{XX}^\top \quad (\text{Lemma}) \\
 & \Rightarrow K\mathbf{C} \frac{1}{T} \mathbf{XX}^\top \succeq \sigma^2(K\mathbf{C} + \sigma^2 \mathbf{I}) \quad (\text{This is true because of eq. 48.}) \\
 & \Rightarrow K\mathbf{C} \frac{1}{T} \mathbf{XX}^\top \succeq \sigma^2 K\mathbf{C} + \sigma^4 \mathbf{I} \\
 & \Rightarrow \frac{1}{K} \mathbf{C}^{-1} \frac{1}{T} \mathbf{XX}^\top \frac{1}{K} \mathbf{C}^{-1} \succeq \sigma^2 \frac{1}{K} \mathbf{C}^{-1} \frac{1}{K} \mathbf{C}^{-1} + \sigma^4 \frac{1}{K} \mathbf{C}^{-1} \frac{1}{K} \mathbf{C}^{-1} \frac{1}{K} \mathbf{C}^{-1} \quad (\text{Lemma}) \\
 & \Rightarrow \mathbf{P}^\top \frac{1}{T} \mathbf{XX}^\top \mathbf{P} \succeq \sigma^2 \mathbf{P}^\top \mathbf{P} + \sigma^4 \mathbf{P}^\top \mathbf{P} \mathbf{P}^\top \mathbf{P} \\
 & \Rightarrow \mathbf{P}^\top \left(\frac{1}{T} \mathbf{XX}^\top + \frac{1}{4} K\mathbf{C} \right) \mathbf{P} \succeq \frac{1}{4} \mathbf{I} + \sigma^2 \mathbf{P}^\top \mathbf{P} + \sigma^4 \mathbf{P}^\top \mathbf{P} \mathbf{P}^\top \mathbf{P} \\
 & \Rightarrow \mathbf{U} \mathbf{S}^4 \mathbf{U}^\top \succeq \frac{1}{4} \mathbf{I} + \sigma^2 \mathbf{P}^\top \mathbf{P} + \sigma^4 \mathbf{P}^\top \mathbf{P} \mathbf{P}^\top \mathbf{P} \\
 & \Rightarrow \mathbf{S}^4 \succeq \left(\frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U} \right) \left(\frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U} \right) \\
 & \Rightarrow (\mathbf{S}^2 - \left(\frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U} \right)) \mathbf{S}^2 + \left(\frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U} \right) (\mathbf{S}^2 - \left(\frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U} \right)) \succeq \mathbf{0}.
 \end{aligned}$$

Denoting $\mathbf{A} = \frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U}$ and $\mathbf{B} = \mathbf{S}^2 - (\frac{1}{2} \mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U})$, we can simplify the above inequality as

$$\mathbf{BS}^2 + \mathbf{AB} \succeq \mathbf{0}. \quad (57)$$

We now prove that if inequality 57 is true, then $\mathbf{B} \succeq \mathbf{0}$, using proof by contradiction:

If $\mathbf{B} \not\succeq \mathbf{0}$, then there exists a eigenvector \mathbf{x} s.t. $\mathbf{Bx} = \lambda \mathbf{x}$, $\lambda < 0$. Then $\mathbf{x}^\top \mathbf{BS}^2 \mathbf{x} = \lambda \mathbf{x}^\top \mathbf{S}^2 \mathbf{x} \leq 0$ due to the p.s.d. of \mathbf{S}^2 . Similarly, for the same eigenvector \mathbf{x} , we could arrive at the same conclusion for \mathbf{AB} ,

i.e., $\mathbf{x}^\top \mathbf{A}\mathbf{B}\mathbf{x} = \lambda \mathbf{x}^\top \mathbf{A}\mathbf{x} \leq 0$. This implies that $\mathbf{x}^\top (\mathbf{B}\mathbf{S}^2 + \mathbf{A}\mathbf{B})\mathbf{x} \leq 0$ which is contradict to inequality 57. Therefore $\mathbf{B} \succeq 0$. Now we could continue the deduction as follows,

$$\begin{aligned}\mathbf{B} &\succeq 0 \\ \Rightarrow \mathbf{S}^2 &\succeq \frac{1}{2}\mathbf{I} + \sigma^2 \mathbf{U}^\top \mathbf{P}^\top \mathbf{P} \mathbf{U} \\ \Rightarrow \mathbf{U} \mathbf{S}^2 \mathbf{U}^\top &\succeq \frac{1}{2}\mathbf{I} + \sigma^2 \mathbf{P}^\top \mathbf{P} \\ \Rightarrow \mathbf{P}^{-\top} \mathbf{U} \mathbf{S}^2 \mathbf{U}^\top \mathbf{P}^{-1} &\succeq \frac{1}{2}K\mathbf{C} + \sigma^2 \mathbf{I} \\ \Rightarrow \mathbf{Q} - \sigma^2 \mathbf{I} &\succeq 0.\end{aligned}$$

796

□

797 We can easily ensure that $K\mathbf{C} - \sigma^2 \mathbf{I}$ is p.s.d., which guarantees that $\mathbf{Q} - \sigma^2 \mathbf{I}$ is p.s.d. according to the
798 **Theorem**. Therefore $\mathbf{L} = \mathbf{W}_{:,K} \mathbf{D}_{:K,:K}^{\frac{1}{2}}$ is valid.

799 **Code Availability**

800 The code is available at <https://github.com/waq1129/brainkernel>.

801 **Data Availability**

802 The HCP datasets are publicly available at <https://www.humanconnectome.org/>. The visual cortex
803 dataset is publicly available at https://nilearn.github.io/modules/generated/nilearn.datasets.fetch_haxby.html. The BOLD5000 dataset is publicly available at <https://bold5000.github.io/>.
804 The Sherlock movie dataset is publicly available at <https://dataspace.princeton.edu/handle/88435/dsp01nz8062179>.
805
806

807 **Acknowledgments**

808 This work was supported by the McKnight Foundation (JP), NSF CAREER Award IIS-1150186 (JP), the
809 Simons Collaboration on the Global Brain (SCGB AWD1004351) (JP) and a J. Insley Blair Pyne Fund
810 Award (to JP, BE, KN).

References

1. Schäfer J, Strimmer K, et al. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical applications in genetics and molecular biology*. 2005;4(1):32.
2. Bickel PJ, Levina E. Regularized estimation of large covariance matrices. *The Annals of Statistics*. 2008; p. 199–227.

3. Hsieh CJ, Sustik MA, Dhillon IS, Ravikumar PK, Poldrack R. BIG & QUIC: Sparse inverse covariance estimation for a million variables. In: NIPS; 2013. p. 3165–3173.
4. Treister E, Turek JS. A block-coordinate descent approach for large-scale sparse inverse covariance estimation. In: NIPS; 2014. p. 927–935.
5. Varoquaux G, Gramfort A, Poline JB, Thirion B. Brain covariance selection: better individual functional connectivity models using population prior. In: NIPS; 2010. p. 2334–2342.
6. Van Essen DC, Smith SM, Barch DM, Behrens TE, Yacoub E, Ugurbil K, et al. The WU-Minn human connectome project: an overview. *Neuroimage*. 2013;80:62–79.
7. Brett M, Johnsrude IS, Owen AM. The problem of functional localization in the human brain. *Nature reviews neuroscience*. 2002;3(3):243.
8. Stein ML. Interpolation of spatial data: some theory for kriging. Springer Science & Business Media; 2012.
9. Rasmussen CE. Gaussian processes in machine learning. In: Advanced lectures on machine learning. Springer; 2004. p. 63–71.
10. Lawrence ND. Gaussian process latent variable models for visualisation of high dimensional data. In: Advances in neural information processing systems; 2004. p. 329–336.
11. Damianou A, Lawrence N. Deep gaussian processes. In: Artificial Intelligence and Statistics; 2013. p. 207–215.
12. Kitterle FL, Kaye RS. Hemispheric symmetry in contrast and orientation sensitivity. *Attention, Perception, & Psychophysics*. 1985;37(5):391–396.
13. Di Lollo V. Hemispheric symmetry in duration of visible persistence. *Perception & Psychophysics*. 1981;29(1):21–25.
14. Westcott M. Hemispheric symmetry of the EEG during the Transcendental Meditation technique. Department of Psychology, University of Durham, Durham, England. 1973;
15. Fan J, Fan Y, Lv J. High dimensional covariance matrix estimation using a factor model. *Journal of Econometrics*. 2008;147(1):186–197.
16. Fan J, Liao Y, Liu H. An overview of the estimation of large covariance and precision matrices. *The Econometrics Journal*. 2016;19(1):C1–C32.
17. Lawrence ND. Learning for Larger Datasets with the Gaussian Process Latent Variable Model. In: AISTATS. vol. 11; 2007. p. 243–250.
18. Damianou AC, Titsias MK, Lawrence ND. Variational inference for uncertainty on the inputs of gaussian process models. arXiv preprint arXiv:14092287. 2014;.
19. Hensman J, Fusi N, Lawrence ND. Gaussian processes for big data. arXiv preprint arXiv:13096835. 2013;.
20. Wu TT, Lange K. Coordinate descent algorithms for lasso penalized regression. *The Annals of Applied Statistics*. 2008; p. 224–244.
21. Wang H. Coordinate descent algorithm for covariance graphical lasso. *Statistics and Computing*. 2014;24(4):521–529.
22. Tseng P, Yun S. A coordinate gradient descent method for nonsmooth separable minimization. *Mathematical Programming*. 2009;117(1):387–423.

23. Fox MD, Snyder AZ, Zacks JM, Raichle ME. Coherent spontaneous activity accounts for trial-to-trial variability in human evoked brain responses. *Nature neuroscience*. 2006;9(1):23–25.
24. Smith SM, Fox PT, Miller KL, Glahn DC, Fox PM, Mackay CE, et al. Correspondence of the brain's functional architecture during activation and rest. *Proceedings of the National Academy of Sciences*. 2009;106(31):13040–13045.
25. Cole MW, Ito T, Bassett DS, Schultz DH. Activity flow over resting-state networks shapes cognitive task activations. *Nature Neuroscience*. 2016;.
26. WU-Minn H. 1200 subjects data release reference manual; 2017.
27. Delgado MR, Nystrom LE, Fissell C, Noll D, Fiez JA. Tracking the hemodynamic responses to reward and punishment in the striatum. *Journal of neurophysiology*. 2000;84(6):3072–3077.
28. Binder JR, Gross WL, Allendorfer JB, Bonilha L, Chapin J, Edwards JC, et al. Mapping anterior temporal lobe language areas with fMRI: a multicenter normative study. *Neuroimage*. 2011;54(2):1465–1475.
29. Buckner RL, Krienen FM, Castellanos A, Diaz JC, Yeo BT. The organization of the human cerebellum estimated by intrinsic functional connectivity. *Journal of neurophysiology*. 2011;.
30. Hariri AR, Brown SM, Williamson DE, Flory JD, De Wit H, Manuck SB. Preference for immediate over delayed rewards is associated with magnitude of ventral striatal activity. *Journal of Neuroscience*. 2006;26(51):13213–13217.
31. Smith R, Keramati K, Christoff K. Localizing the rostral lateral prefrontal cortex at the individual level. *Neuroimage*. 2007;36(4):1387–1396.
32. Barch DM, Burgess GC, Harms MP, Petersen SE, Schlaggar BL, Corbetta M, et al. Function in the human connectome: task-fMRI and individual differences in behavior. *Neuroimage*. 2013;80:169–189.
33. Haxby JV, Gobbini MI, Furey ML, Ishai A, Schouten JL, Pietrini P. Distributed and overlapping representations of faces and objects in ventral temporal cortex. *Science*. 2001;293(5539):2425–2430.
34. Esteban O, Blair R, Markiewicz CJ, Berleant SL, Moodie C, Ma F, et al.. poldracklab/fmriprep: 1.0.0-rc5; 2017. Available from: <https://doi.org/10.5281/zenodo.996169>.
35. Abraham A, Pedregosa F, Eickenberg M, Gervais P, Mueller A, Kossaifi J, et al. Machine learning for neuroimaging with scikit-learn. *Frontiers in neuroinformatics*. 2014;8:14.
36. Chen J, Leong YC, Honey CJ, Yong CH, Norman KA, Hasson U. Shared memories reveal shared structure in neural activity across individuals. *Nature neuroscience*. 2017;20(1):115.
37. Simony E, Honey CJ, Chen J, Lositsky O, Yeshurun Y, Wiesel A, et al. Dynamic reconfiguration of the default mode network during narrative comprehension. *Nature communications*. 2016;7:12141.
38. Gershman SJ, Blei DM, Pereira F, Norman KA. A topographic latent source model for fMRI data. *NeuroImage*. 2011;57(1):89–100.
39. Manning JR, Ranganath R, Norman KA, Blei DM. Topographic factor analysis: a bayesian model for inferring brain networks from neural data. *PloS one*. 2014;9(5):e94914.
40. Wu A, Pashkovski S, Datta SR, Pillow JW. Learning a latent manifold of odor representations from neural responses in piriform cortex. In: *Advances in Neural Information Processing Systems*; 2018. p. 5378–5388.

41. Chang N, Pyles JA, Marcus A, Gupta A, Tarr MJ, Aminoff EM. BOLD5000, a public fMRI dataset while viewing 5000 visual images. *Scientific data*. 2019;6(1):49.
42. Grosenick L, Klingenberg B, Katovich K, Knutson B, Taylor JE. Interpretable whole-brain prediction analysis with GraphNet. *NeuroImage*. 2013;72:304–321.
43. Rao N, Cox C, Nowak R, Rogers TT. Sparse overlapping sets lasso for multitask learning and its application to fmri analysis. *Advances in neural information processing systems*. 2013;26:2202–2210.
44. Wu A, Park M, Koyejo OO, Pillow JW. Sparse Bayesian structure learning with dependent relevance determination priors. In: Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, editors. *Advances in Neural Information Processing Systems 27*. Curran Associates, Inc.; 2014. p. 1628–1636. Available from: <http://papers.nips.cc/paper/5233-sparse-bayesian-structure-learning-with-dependent-relevance-determination-priors.pdf>.
45. Wu A, Koyejo O, Pillow J. Dependent relevance determination for smooth and structured sparse regression. *Journal of Machine Learning Research*. 2019;20(89):1–43.
46. Xu H, Lorbert A, Ramadge PJ, Guntupalli JS, Haxby JV. Regularized hyperalignment of multi-set fMRI data. In: 2012 IEEE Statistical Signal Processing Workshop (SSP). IEEE; 2012. p. 229–232.
47. Haxby JV, Guntupalli JS, Nastase SA, Feilong M. Hyperalignment: Modeling shared information encoded in idiosyncratic cortical topographies. *ELife*. 2020;9:e56601.
48. Lawrence ND. Gaussian Process Latent Variable Models for Visualisation of High Dimensional Data. In: Nips. vol. 2; 2003. p. 5.
49. Nocedal J. Updating quasi-Newton matrices with limited storage. *Mathematics of computation*. 1980;35(151):773–782.
50. Belkin M, Niyogi P. Laplacian eigenmaps for dimensionality reduction and data representation. *Neural computation*. 2003;15(6):1373–1396.
51. John J, Draper N. An alternative family of transformations. *Applied Statistics*. 1980; p. 190–197.
52. Stein ML. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer; 1999.
53. Lázaro-Gredilla M, Quiñonero-Candela J, Rasmussen CE, Figueiras-Vidal AR. Sparse spectrum Gaussian process regression. *The Journal of Machine Learning Research*. 2010;11:1865–1881.
54. Wu A, Aoi MC, Pillow JW. Exploiting gradients and Hessians in Bayesian optimization and Bayesian quadrature. *arXiv preprint arXiv:170400060*. 2017;.
55. Tipping ME, Bishop CM. Probabilistic principal component analysis. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 1999;61(3):611–622.
56. Vodrahalli K, Chen PH, Liang Y, Baldassano C, Chen J, Yong E, et al. Mapping between fMRI responses to movies and their natural language annotations. *NeuroImage*. 2018;180:223–231.

Appendices

A Relationship between brain kernel and deep GPs

The BK model can be viewed as a sequence of two nonlinear functions, each with a Gaussian process prior.

In the main text, we explained that the latent embedding \mathbf{z} is obtained by applying a function f over the input 3D coordinate \mathbf{x} , where $f \sim \mathcal{GP}(m, k)$. This is the first layer GP. The second component of the brain kernel model is a probability distribution over neural activity as a function of location in embedding space. Our modeling assumption is that neural activity changes smoothly as a function of location in embedding space, or equivalently, that correlations in neural activity fall off smoothly with distance in latent space. This motivates the use of a second GP to model each vector of neural activity. Thus:

$$g \sim \mathcal{GP}(0, \kappa_g), \quad (58)$$

$$\mathbf{v} = g(\mathbf{z}), \quad (59)$$

where $\kappa_g(\mathbf{z}_i, \mathbf{z}_j) = \exp(-\|\mathbf{z}_i - \mathbf{z}_j\|^2/2)$ is a standard RBF covariance function (marginal variance and length scale = 1), and $\mathbf{v} \in \mathbb{R}^n$ denotes a vector of n voxels' neural activity at latent locations $\{\mathbf{z}_1, \dots, \mathbf{z}_n\}$. Given that latent voxel locations are themselves a nonlinear function of the true voxel locations, this allows to write, equivalently: $\mathbf{v} = g(f(\mathbf{x}_1, \dots, \mathbf{x}_n))$. Therefore, we could eventually obtain the neural activity \mathbf{v} by applying a two-layer GPs to the input 3D coordinates \mathbf{X} .

B More decoding results

B.1 HCP database

In this section, we present the classification accuracy performance for the motor task, the language processing task, the emotion processing task, the relational processing task, and the social cognition task.

Table 1: The basic information of the fMRI scans [26]

Task	Runs	Frames per run	Run Duration (min:sec)
Working Memory	2	405	5:01
Gambling	2	253	3:12
Motor	2	284	3:34
Language Processing	2	316	3:57
Social Cognition	2	274	3:27
Relational Processing	2	232	2:56
Emotion Processing	2	176	2:16

Similar to the working memory and gambling tasks, each task fMRI was from the same HCP project as the resting-state fMRI. They shared the same coordinate system and preprocessing pipeline. The task fMRI was aligned with the MNI template with the same 59,412 voxels as the resting-state fMRI data. Instead of analyzing the whole-brain data for brain decoding, we worked with functional ROIs whose averaged activity levels were above a certain threshold for all tasks.

We formulated each task as a binary classification problem. We used the same Bayesian linear regression classifiers as described in the working memory and gambling sections. We trained the classifier with three priors on one run and calculated the accuracy performance on the second run, then switched the training and test runs. This procedure was repeated ten times. The ± 1 labels were treated as continuous target values during training, and the test accuracy was evaluated by taking the sign of the prediction. We computed the averaged accuracies across two runs and 10 repetitions for the three priors for all subjects in each dataset. Results are presented in Fig. 14 to Fig. 18.

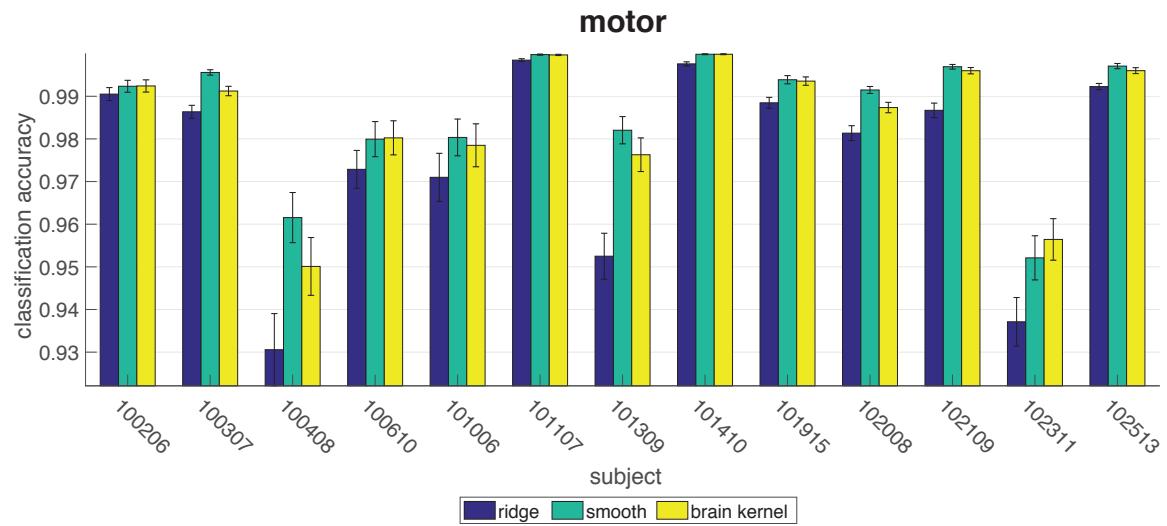


Fig 14. Accuracy performance on the motor task.

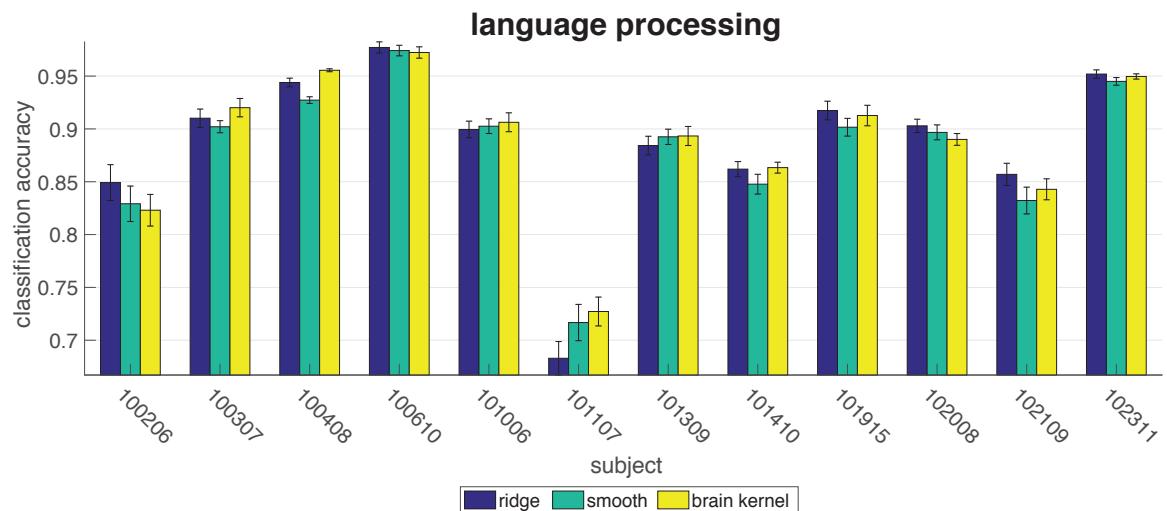


Fig 15. Accuracy performance on the language processing task.

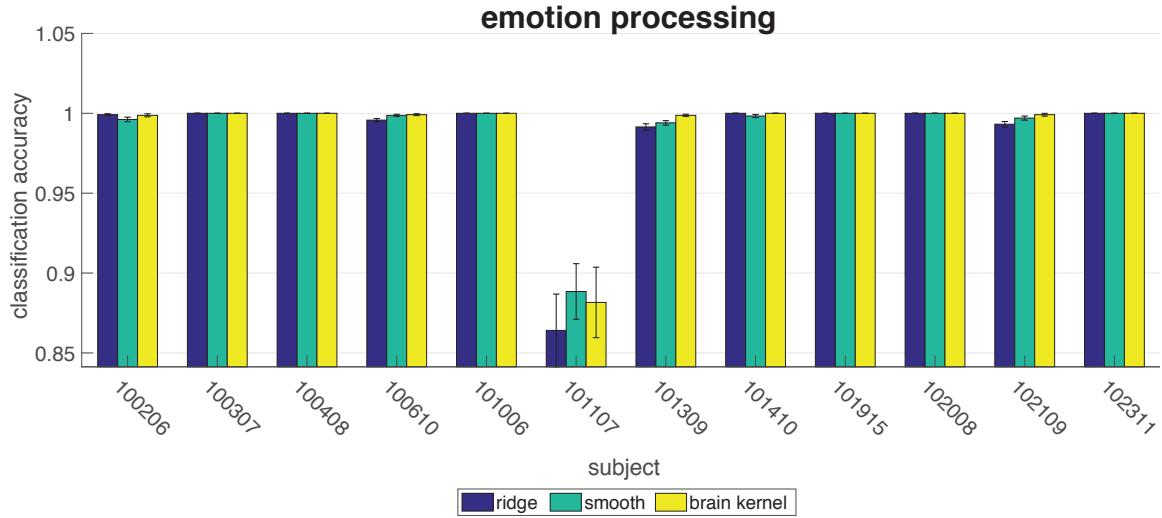


Fig 16. Accuracy performance on the emotion processing task.

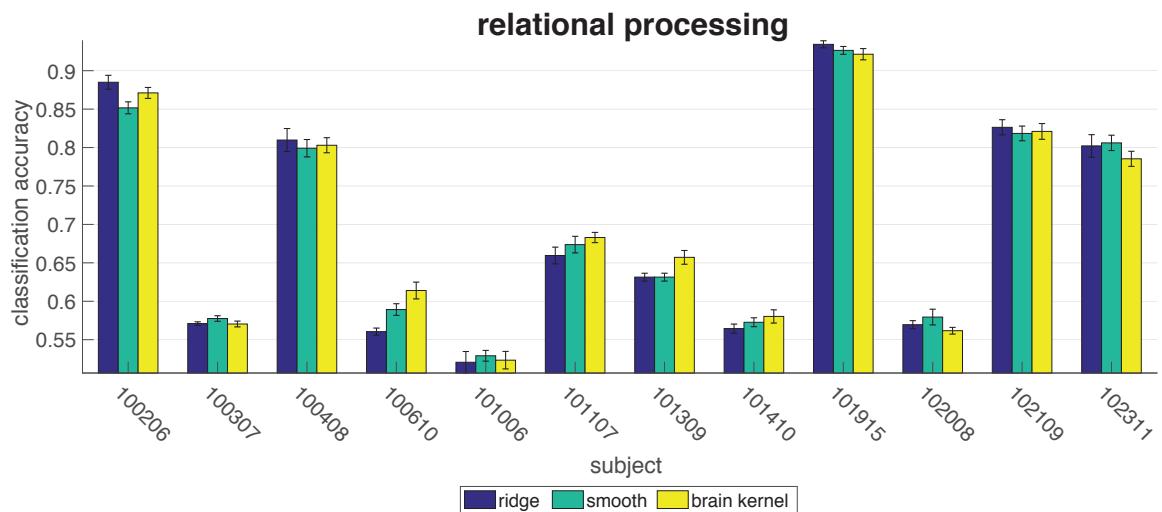


Fig 17. Accuracy performance on the relational processing task.

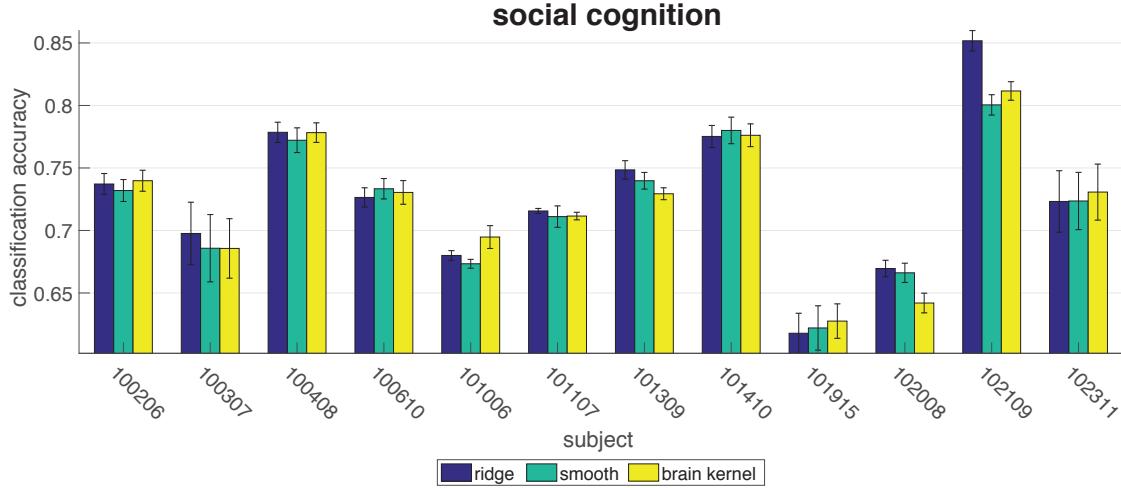


Fig 18. Accuracy performance on the social cognition task.

B.2 Sherlock movie fMRI dataset

We examined the Sherlock fMRI dataset, in which participants were scanned while they watched the British television program “Sherlock” for 50 min [36]. The fMRI data comprised 1,973 TRs (Repetition Time), where each TR was 1.5 s of the movie. Before performing any analysis, the fMRI data were preprocessed and aligned to MNI space using the techniques described in prior work [36]. We examined the brain data averaged across all subjects to smooth out individual variability. We identified 11 ROIs previously implicated in processing naturalistic stimuli, comprising the default mode network (DMN-A, DMN-B), the ventral and dorsal language areas, and the primary auditory and visual cortices [37].

For each ROI, we performed a standard decoding task to relate the fMRI signal to the representation of the semantic content of each movie frame [56]. We trained the decoding model using the first half of the movie and tested on the second half of the movie. The decoding task was called “scene classification”. We divided up the second half of the movie into 25 uniformly-sized chunks, and used the correlation to match predicted semantic content with the ground truth for each chunk, and reported the percentage of the chunks that the match is perfect. Since there are 25 chunks, random chance performance at this task is 4%. Fig. 19 presents the classification accuracy performance. We observed that the smoothing prior and the brain kernel prior both converged to the ridge prior after hyperparameter estimation (3-fold cross validation during training), i.e., the estimated length scale was very small and all three priors had the same performance. Therefore, smoothness did not help the scene classification task.

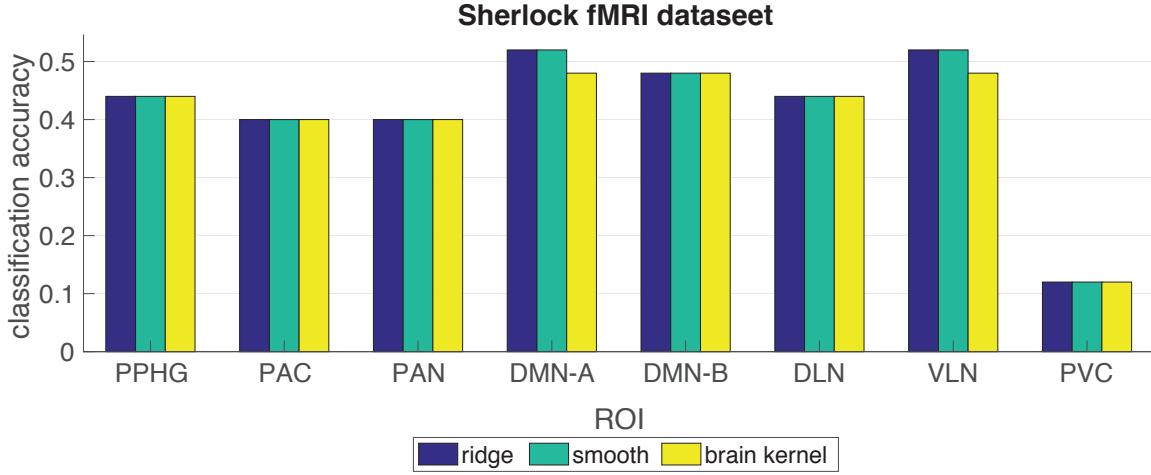


Fig 19. Accuracy performance on the Sherlock movie fMRI dataset. The ROI information is the same as Fig. 12

B.3 BOLD5000 dataset

BOLD5000 is a public fMRI dataset in which subjects viewed 5000 visual images [41]. We used the fMRI datasets collected for 4 subjects with 15 sessions per subject and 333 TRs per session. We extracted predefined functional ROIs [41] from the whole-brain data with 3804 voxels. The decoding task was binary classification. The binary labels were living objects versus nonliving objects. For each subject, we used 14 sessions for training and 1 session for test. Thus the final accuracy was an average over 15 runs by treating each session as the test set once (Fig. 20). We observed that the brain kernel did not reliably outperform the ridge prior and the smoothing prior estimates across subjects.

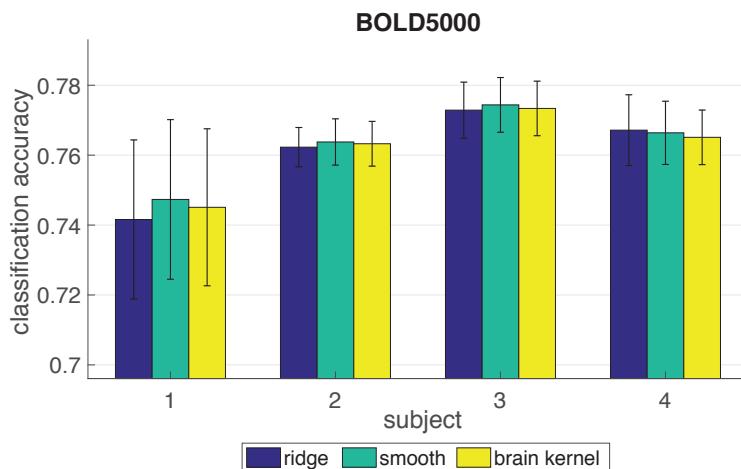


Fig 20. Accuracy performance on the BOLD5000 dataset.

C More factor modeling results

C.1 Visual recognition task

We examined the visual recognition task fMRI dataset with the same Bayesian FA model. We implemented the same “co-smoothing” evaluation as described in the main section. The only difference is, instead of splitting the runs into halves for training and test within each subject, we treated one subject as the test set and the other five subjects as the training set. This is consistent with the inter-subject classification in the decoding task. We report the normalized test R^2 for the three priors when treating each subject as the test set in Fig. 21. For most subjects, the brain kernel didn’t show a dominating performance over both the ridge covariance and the smooth RBF kernel.

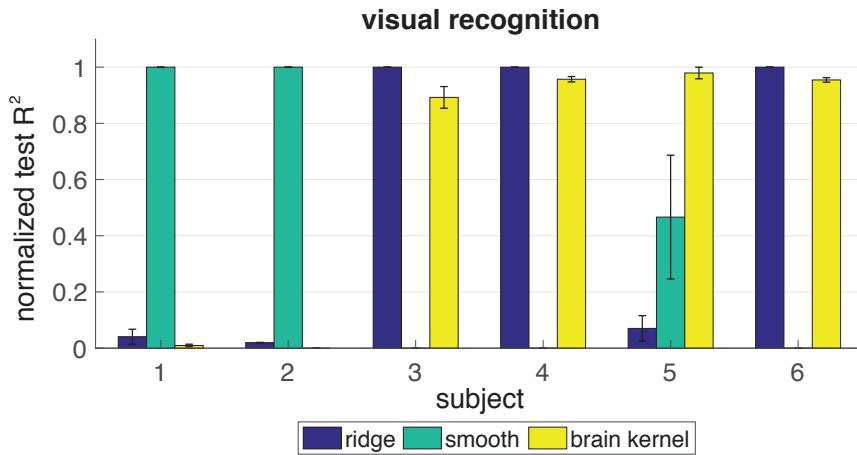


Fig 21. Normalized test R^2 performance on the visual recognition task.

C.2 BOLD5000 dataset

The BOLD5000 dataset consists of 4 subjects with 15 sessions per subject and 333 TRs per session. We extracted predefined functional ROIs [41] from the whole-brain data with 3804 voxels. When evaluating the performance using “co-smoothing”, we split the 15 sessions into halves within each subject, one for training and one for inference and test as described in the main section. We report the normalized test R^2 for the three priors for each subject in Fig. 22. For most subjects, the brain kernel didn’t show a dominating performance over both the ridge covariance and the smooth RBF kernel. The smoothing prior had better performance than the brain kernel which was better than the ridge prior. This result is consistent with the decoding performance in Fig. 20.

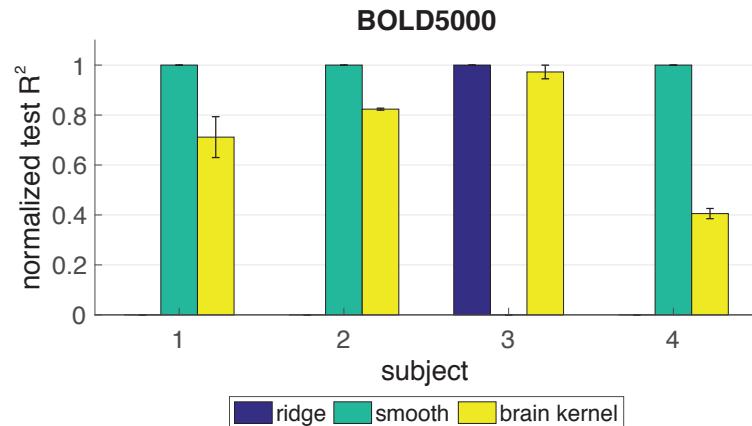


Fig 22. Normalized test R^2 performance on the BOLD5000 dataset.