

Supporting Information

Margulies et al. 10.1073/pnas.1608282113

SI Materials and Methods

Connectivity Data. Connectivity data for the human and macaque monkey brains were assembled from openly available sources. Because acquisition procedures and the method for assessing connectivity varied between species, the first step was to transform each species' dataset into an affinity matrix describing the pairwise similarity of connectivity values.

Human. The connectivity matrix for the human brain was based on 1 h of resting-state fMRI data acquired through the HCP (39) and made publicly available for download on ConnectomeDB (71). In brief, for each individual, a functional connectivity matrix was calculated using the correlation coefficient across four minimally preprocessed (40, 72–74), spatially normalized, and concatenated 15-min resting-state fMRI scans. We began with the group-averaged functional connectivity data, which include 820 individuals coregistered using MSMAll. More information about this dataset can be found at www.humanconnectome.org/documentation/S900/820_Group-average_fmri_Connectivity_December2015.pdf. The “dense” functional connectome matrix consists of pairwise z-transformed correlation values between all grayordinates (91,282 rows). Cortical data are represented in HCP 32k_LR surface space (75), consisting of 32,492 total nodes per hemisphere (59,412 excluding the medial wall). Subcortical structures are represented in volumetric space and included 31,870 voxels.

We began by transforming the z to r correlation values with a hyperbolic tangent function, which scales them between -1 and 1 . For each row in the matrix, the values of the top 10% of connections were retained, whereas all others were zeroed. Remaining connections were almost all positive, except for $\sim 5,000$ (less than $10^{-6}\%$ of all connections) with negative values from 23 voxels. The voxels with negative connections were all located in ventral subcortical regions. These connections were zeroed as well. Because this procedure rendered the connectivity matrix A asymmetric, similarity between all pairs of rows was calculated using cosine distance, resulting in the positive, symmetric affinity matrix L . The full affinity matrix, L , consisted of weights between zero and one, representing similarity of connectivity profiles among over 4 billion node pairs.

Macaque Monkey. Macaque monkey connectivity data were derived from the axonal tract-tracing connectivity database CoCoMac (67, 68) and presented in the F99 template space (76) using cortical area labels based on the Bonin–Bailey parcellation (77, 78). One cortical area was removed, because it displayed no connectivity with other areas, leaving 25 areas remaining. The 25×25 connectivity matrix initially consists of sources (rows) and targets (columns), where edges were values between one (weak) and three (strong) connections. Of the possible edges, 56% were present in the connectivity matrix. To create a similarity matrix that accounted for bidirectional connectivity, the connectivity matrix, A , was transposed and concatenated with the initial matrix: $A = (A, A^T)$. A was then transformed into a positive, symmetric affinity matrix L by calculating the Euclidean distance between each pair of rows.

Connectivity Embedding. We used diffusion embedding (42), a nonlinear dimensionality reduction technique, to recover a low-dimensional embedding from high-dimensional connectivity data. Connectivity data include both local and long-range connections. Diffusion maps translate these relationships into distances and represent the global connectivity structure as a distribution of cortical points in an embedding space. Cortical points that are

strongly connected by either many connections or few very strong connections are close in this space, whereas points without connections are far apart. Linear techniques, such as principal component analysis, are unable to project such data without appropriate kernel manipulations. Among the numerous nonlinear dimensionality reduction algorithms currently in use, we chose diffusion map embedding (42), because the diffusion process limits the distances of influence to the graph neighborhood, thereby ensuring a stable representation of connections, regardless of the graph size (79). Other similar approaches have been previously applied to whole-brain structural connectivity data (80, 81).

Diffusion maps are a one-parameter (α) family of graph Laplacians that integrate local information into a global description. A diffusion map embedding reduces high-dimensional data to a low-dimensional representation that combines geometry with probability distribution of data points. Using an appropriately normalized random walk process, the parameter, α , can control whether the low-dimensional embedding reflects the geometry of the set, regardless of density of points or the long-time dynamics of the samples without uniform sampling. In other words, the parameter α controls the influence of density of sampling points on the underlying manifold ($\alpha = 0$, maximal influence of sampling density; $\alpha = 1$, no influence of sampling density). Their relationship to other embedding methods based on graph Laplacians is detailed in ref. 82.

Diffusion maps, which use an α of 0.5, are well-suited for the analysis of brain connectivity data for multiple reasons. First, they retain the global relations between data points in the embedded space. Second, they are more robust to noise in the connectivity matrix, unlike other techniques, such as Isomap (83). Third, by using the appropriate choice of α , they can be made less sensitive to the distribution of data. Fourth, decreasing eigenvalues reflect a natural ordering of the diffusion process; the largest eigenvalues correspond to the slowest processes and thereby, represent the slowest variance in connectivity patterns. Fifth, by relying on local distances, they address the curse of dimensionality, because smaller distances are more meaningful than larger distances as the number of dimensions increases.

Technically, following ref. 42, the species-specific connectivity matrices (represented as L) described previously can be transformed to define a Markov chain with transition matrix P . The eigenvectors of this transition matrix define an embedding that results in a representation of each vertex or region as a point in the embedding space. The mutual Euclidean distances between these representations reflect how closely connected the corresponding vertices are in the graph with regard to the diffusion process defined by the transition matrix. See algorithm 1 for pseudocode. Increasing the diffusion time (t) allows examination of the intrinsic geometric structure of the data at larger and larger scales, allowing a balance between details in the input data and the scale of assessment (79). In this implementation, we used an automated estimation of the diffusion time using a damped regularization process. In this approach, the eigenvalues λ_i are divided by $1 - \lambda_i$ to provide robustness against small, noisy eigenvalues. Such diffusion map embedding has been previously applied to functional task and resting MRI data (84–89). The source code for this method is available at <https://github.com/satra/mapalign>.

Geodesic Distance Along the Cortical Surface. Geodesic distance along the cortical surface was calculated using an algorithm that approximates the exact distance along the shortest path between

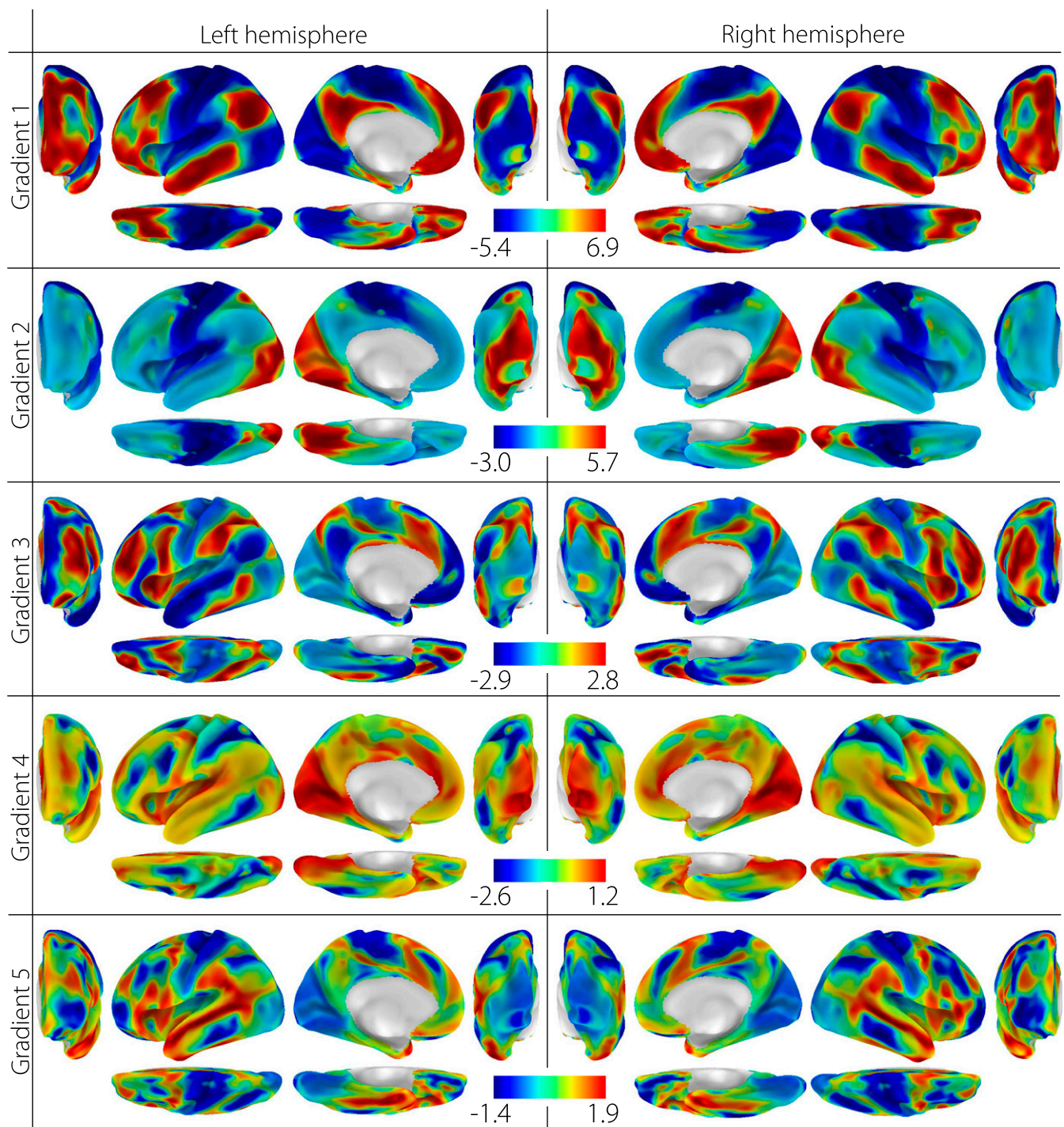


Fig. S1. Human connectivity gradients 1–5. The first five components result from diffusion embedding of the human connectivity matrix. The first five are shown because of the drop in variance explained after the fifth component (Fig. S2).

Behavioral Domains (BD)

Brainmap behavior domains

Percentile along gradient

likelihood-ratio (using z-stat threshold > 2.3)

2.6

1.1

5 of 6

www.pnas.org/cgi/content/short/1608282113