# Functional Neurodata Graph Service: A One-Click Pipeline and Web Service for Reliable Functional Connectome Estimation

Eric W. Bridgeford[1,2,*], Gregory Kiar[2,3], Tanay Agarwal[1], Eric Walker[1,2], and Joshua T. Vogelstein[2]

[1]Department of Computer Science, Johns Hopkins University

[2]Department of Biomedical Engineering, Johns Hopkins University

[3]Department of Biomedical Engineering, McGill University

[*]ericwb95@gmail.com

## 1  Abstract

In recent years, functional magnetic resonance imaging (fMRI) analysis has become one of the most widely used techniques for assessing cognitive function. fMRI allows an unparalleled combination of temporal and spatial acuity, providing researchers a glance at the dynamic activity of individual regions in the brain. Unfortunately, computational restrictions of collecting and analyzing an fMRI dataset necessitate technological investments and expertise for interested researchers. Through the Functional NeuroData Graph Service, we develop a robust fMRI processing pipeline for providing automated acquisition of functional connectomes. Furthermore, we provide an open-source scientific container pre-loaded with our software and all dependencies to facilitate the distribution of our methodologies. Finally, we provide utilities and a web-service for deploying our pipeline using Electronic Cloud Computing, enabling researchers to process entire datasets in just a few hours with no external dependencies.

## 2  Introduction

For the past thirty years, fMRI has presented one of the most intriguing challenges for computational neuroscience. While researchers have long understood that an increase in blood flow was necessary for increased brain function (such as during periods of high brain activity), it was not until a seminal study by [1] that researchers were able to image the flow of blood in the brain. Oxygen-rich blood (high in Hb, the oxygenated form of hemoglobin) is slightly more magnetic than deoxygenated blood (high in dHb, or deoxygenated hemoglobin). As Hb has different magnetic properties than dHb, a high strength magnetic

field such as that of an MRI scanner is able to pick up time-dependent fluctuations in the concentration of Hb. While regions of the brain showing high activity see a brief drop in the concentration of oxygen in vessels supplying the region, the arteries rapidly overcompensate by the hemodynamic response [2], which sees a significant increase in the concentration of oxygenated blood pushed towards the regions with high brain activity. This fluctuation in the Blood Oxygenation Level Dependent (BOLD) signal is what is imaged in an fMRI scanner.

Recently, researchers have been incresingly concerned with developing maps of the brain, or connectomes [3]. In the functional sense, the connectome represents the strength, or correlation, with which two particular regions of the brain fire together. We record this strength as the "connectivity" of two regions. Repeating over all regions of the brain, we obtain a map, or network, of the brain. We can then apply network theoretical approaches to gain insight into spatial patterns that arise at the region level [4].

# 3   Methods

## 3.1   Pipeline

The question of identifying an optimal processing pipeline is, unfortunately, heavily non-convex in that the parameter settings of components interact such that testing individual processing steps standalone is unlikely to lead to the optimal processing pipeline [5]. To begin to answer this question, we conducted a full parameter sweep of many of the most popular techniques for fMRI analysis in [6]. We then refined each step using a combination of quantative and qualitative metrics to yield a pipeline which ensures maximal robustness.

### 3.1.1   Preprocessing

Subject-specific artifacts caused by head motion and limitations of the fMRI sensor are known to be confounding factors in fMRI analyses [7]. The workflow in Figure (4) shows the solution the FNGS pipeline uses to correct for spatial artifacts inherent in the fMRI scanning technique.

**Slice Timing**   To collect an individual 4D EPI sequence, a 3D volume is constructed as a combination of individual 2D slices. The 2D slices are collected incrementally; that is, we collect each 2D slice for approximately 10 milliseconds, and the entire 3D volume is complete in about 30 2D slices. This gives a repetition time, or TR, for the volume on the order of 1 to 3 seconds depending on the scanner. In response to a stimulus, we expect a typical brain response on the order of about 16 seconds. This means that during the course of a single volume being collected, we may have a different amount of BOLD present in the first volume than when the last slice is actually collected [9]. While intuitively we might want our observations to represent a "snapshot" of the data at a given time point, instead we have a "sliding snapshot" over the course of one TR; that is, our observations are not all at a fixed point in time.

To begin, we accept the user-provided acquisition sequence of a single 3D volume. Given the acquisi-

(a) fMRI Preprocessing Pipeline
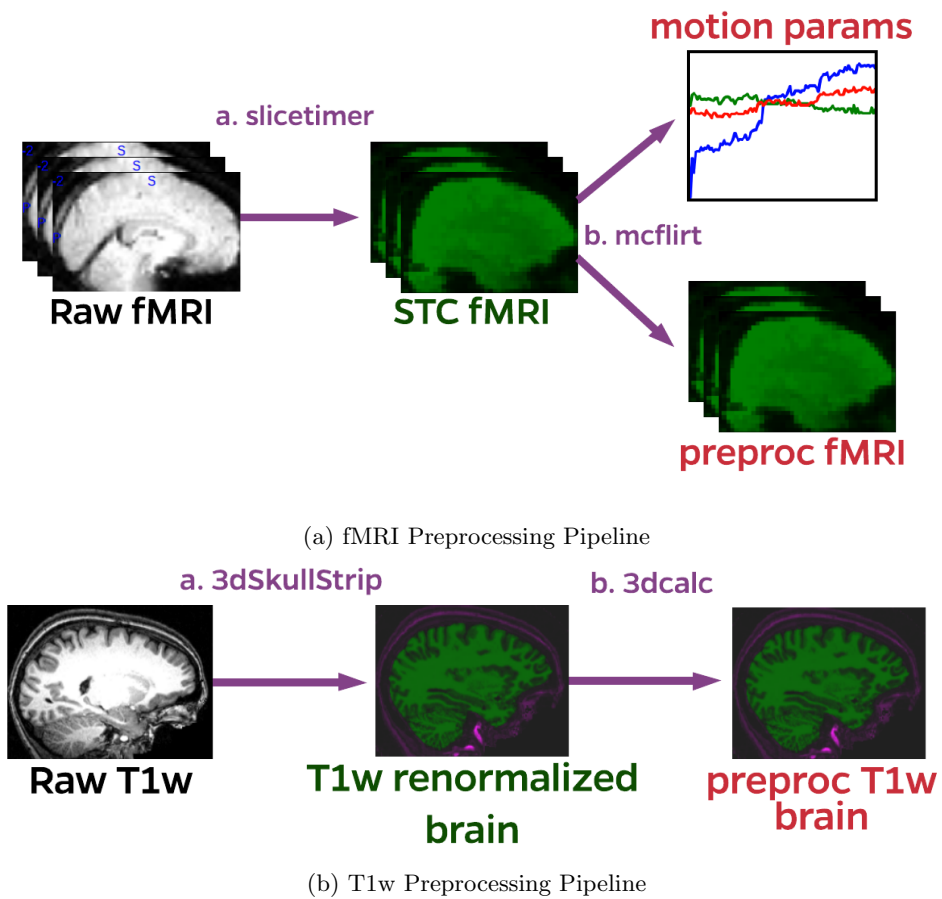


(b) T1w Preprocessing Pipeline

Figure 1: The preprocessing workflow for the FNGS pipeline. Inputs are shown in black, intermediates in green, and results in red, with purple showing computations being performed on our data.

(a) Raw fMRI

(b) Interleaved Acquisition Sequence

(c) Slice-Timing Corrected fMRI
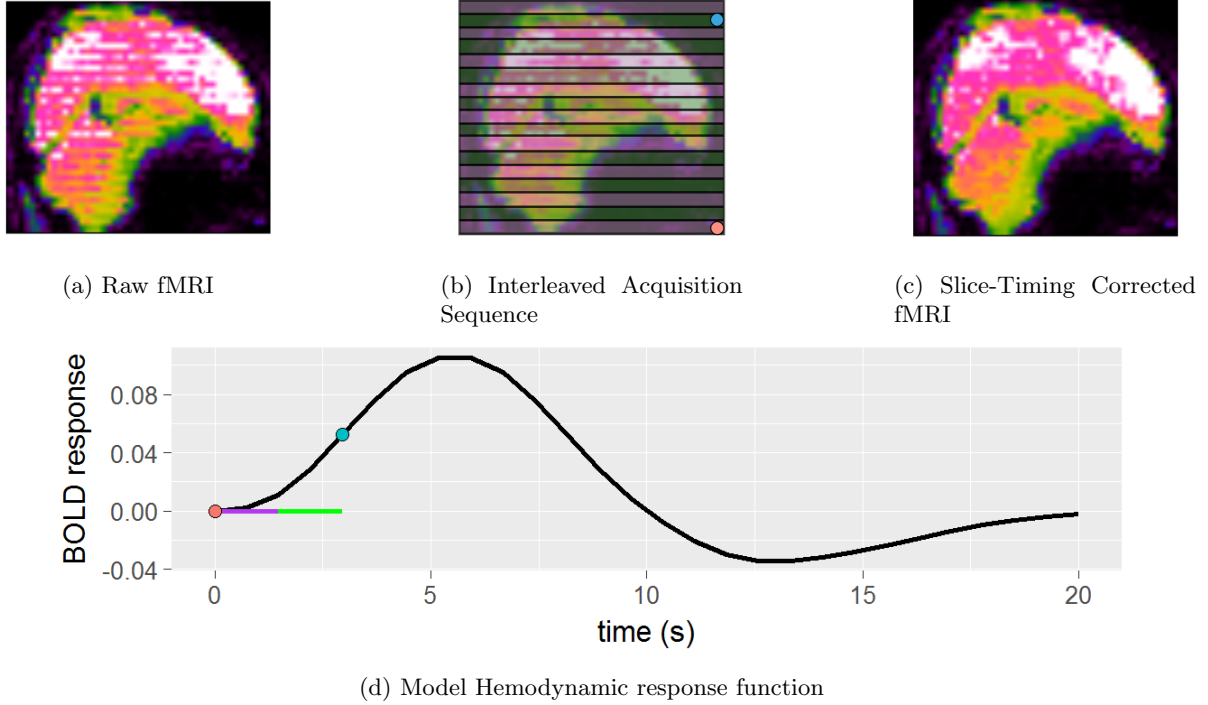
(d) Model Hemodynamic response function

Figure 2: In the raw fMRI shown in (a), we can see the slice-timing effect. Notice the horizontal banding, characteristic of the interleaved-acquisition pattern. In (b), we can see how the interleaved acquisition sequence works. The 3D volume is acquired from bottom to top; first the purple 2D slices, and next the green 2D slices. In (d), we can see an example hemodynamic response, where it is clear that the first slice taken (the red dot) is at a noticeably different time point than the last slice taken (the blue dot) for the first TR (the purple corresponds to the time period when the purple 2D slices are taken, and vice-versa for the green). In (c), we can see the fMRI after slice-timing correction. Data shown is from the HNU1 dataset [8].

tion sequence of the 2D slices, we can compute the TR shift of each 2D slice given the TR information from the header of the brain image. A slice that occurs first in a TR will have a shift of 0, while a slice that occurs at the end of a TR will have a shift of 1. A slice that occurs exactly in the middle of a TR has a shift of 0.5. For each voxel in an individual slice, we use interpolation to re-center our observations to all have a TR shift of 0.5. For details, see Figure (2). We accomplish slicetiming using the slicetimer utility provided by FSL [9]. Optionally, users may choose to skip this step, as slice timing correction is only needed when users are focusing on temporal characteristics of an fMRI session.

**Motion Correction**    During an fMRI session, participants sit in a small, cramped scanner often for between 5 and 10 minutes. During the course of a study, it is fairly common for participants to shift and fidget, even if only small amounts. Small shifts will lead to a person's head being in different spatial positions at each timestep, which will hamper our efforts to standardize the spatial properties of each subject's brain down the line through registration. This is because registration entails aligning each 3D volume in a 4D sequence using the same linear or non-linear transformation. This means that if each 3D volume is not already the same spatially, we may get minor inconsistencies in our registration that will generally decrease our functional connectome quality [10].

Fortunately, given that the subject's brain is the same scaling for each 3D volume (the brain shape itself is not changing in time), we can estimate a 6 degree of freedom rigid affine transformation for each 3D volume (we have 1 translational and 1 rotational parameter per $x, y, z$ direction the subject could move his/her head) using the mean fMRI slice as our reference. Motion correction is implemented using the mcflirt utility [11], which is a simplification of FSL's FLIRT registration tool.

**Anatomical Preprocessing**    To preprocess our anatomical image, we are chiefly concerned with ensuring that we obtain a high quality brain extraction, as many registration techniques will fail with improperly extracted brains. We leverage AFNI's 3dSkullstrip, which provides modifications to the BET algorithm to make it more robust without hyperparameters. Note that 3dSkullstrip renormalizes intensities, so to regain the original intensities, we feed this result as a step function (essentially making it a mask) through 3dcalc and multiply voxelwise with the original image, giving us the original image intensities of the brain and excluding the regions determined to be skull.

### 3.1.2  Registration

Unlike many modalities, fMRI offers a combination of both spatial and temporal acuity. The questions we seek to answer with fMRI generally center around the dynamic activity of each person's brain; we want to see how people's brain activity is similar (or different), free from structure or shape differences. For each subject, however, we have unique structural properties associated with the shape and layout of that particular brain. Using linear and nonlinear realignments, we can reshape the brains of each subject
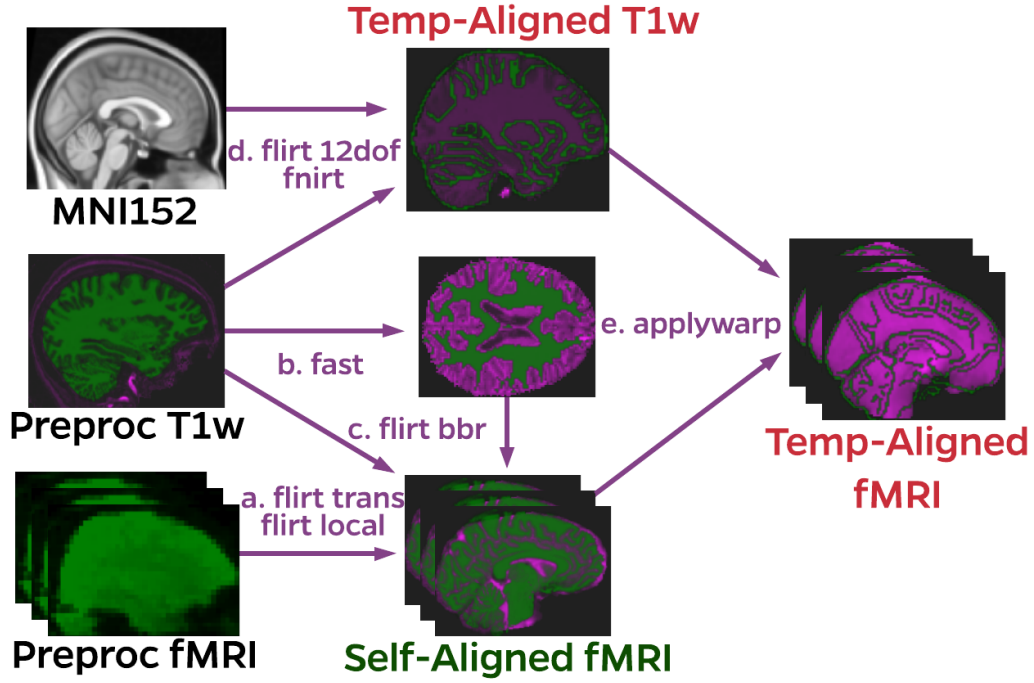
Figure 3: The registration workflow for the FNGS pipeline. We take as input to the pipeline a preprocessed brain, an anatomical brain, and a template brain. We estimate a linear transformation from the fMRI to the structural image to account for the higher resolution of the anatomical image, and estimate a non-linear transformation from the anatomical image to the template image. Finally, we apply the non-linear transformation to the functional brain in the anatomical brain space using our previously estimated transformations for reference.

to that of an average human (or template) brain. With all of the brains in a common spatial domain, we can simply identify the voxels associated with a particular region in the template brain, and then subset these same voxels for all subjects in the dataset. In Figure (3), we can see the registration workflow for the FNGS pipeline.

**Registration Score**   Registration is one of the easiest places for automated fMRI processing to fail. A slight misalignment at one step can have dramatically negative consequences on resulting downstream inferences; we occasionally introduce tearing, distortions, or more serious problems into our scans. To define a quality registration, we introduce the Jaccard index $r$, defined as:

$$r_\gamma(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

for a brain $A$ being aligned to a reference $B$ using registration procedure $\gamma$, where a higher quality registration has a value of $r_\gamma(A, B) > c$ for some cutoff $c$. This statistic captures the intuition of registration in its most basic form; we expect that properly registered brains will have a spatial overlap that is comparatively sized to their spatial union. Often, techniques that are less robust provide better alignments on high-quality data. Since registration failures are not immediately apparent in downstream derivatives, recognizing when one particular technique fails is of the utmost importance.

**Self Registration**   When collecting an fMRI sequence, researchers also collect a higher resolution anatomical scan, using a procedure known as T1-weighting [12]. Using the T1w imaging scheme, researchers can collect imaging data at resolutions exceeding $1\text{mm} \times 1\text{mm} \times 1\text{mm}$, whereas fMRI sequences will often be of substantially lower resolution (generally around $3\text{mm} \times 3\text{mm} \times 3\text{mm}$). This enhanced resolution makes identifying anatomical landmarks far easier using the T1w than the fMRI sequence, which will greatly improve our template registration down the line. To leverage the high quality registration we can obtain from the T1w to the template, we first must know how to reshape the fMRI scan into the same shape of the T1w scan. While the two brains share the same true shape, the T1w and fMRI scans have vastly different dimensions and resolutions that must be accounted for.

To register our input fMRI to our reference T1w image, we begin by first estimating a 3 degree of freedom (DOF) affine transformation with x, y, and z translational parameters with FSL's FLIRT [13] using the *3dtrans.sch* file provided as part of the FSL package. This centers our T1w brain optimally, which will increase our ability to make higher quality registrations donw the line since it will give us a better starting point for later registrations and make us less likely to become stuck in a poor-fit local optima. Next, using this translational transform as our starting point, we estimate a locally-optimized

transformation from our fMRI space to our T1w space. Again, this transformation is heavily robust, and has its hyperparameters tuned to focus on local features of the input (fMRI) and reference (T1w) spaces. We use this cost function because often, input fMRI may have narrow fields of view, resolution constraints, or tearing that will perform poorly using a more global alignment.

Next, we compute a third alignment using our locally-aligned fMRI to our structural T1w space using the bbr cost function provided by FSL [14]. In functional data, the white-matter/gray-matter border is fairly apparent as the gray-matter generally shows higher intensities than the white-matter regions. Leveraging this observation, we can align this white-matter boundary between the fMRI and the T1w scan much better than any anatomical features that may be present in our scans. We first estimate a 6 DOF transformation from our fMRI to our T1w scan, and then segment our T1w scan into a white-matter mask. We then use FLIRT with the bbr (boundary-based registration) cost-function to align the boundaries of the white-matter in the fMRI and T1w scans optimally. This provides a high-quality alignment for intra-modal registration from an EPI to a T1w image.

**Template Registration**   Now that our fMRI sequence is in T1w space, we can align our brain to a reference template. A template brain represents the anatomical average brain of the sampling of subjects it is collected over. This anatomically average brain theoretically represents the average brain we will find during our external investigations, allowing minimal alignment on average. For FNGS, we assume that users will be using the MNI152 template [15].

We first perform a gentle linear transformation of our high-resolution T1w brain using the local-optimisation schedule file from before. Next, we use this local-optimisation registration as the starting point for a more extensive 12-DOF global FLIRT alignment than for the self-registration case. Given that our template brain will theoretically be less similar than simple translations, rotations and scalings can provide, we then use a non-linear registration from our T1w to our template space. We accomplish this using FSL's FNIRT algorithm [16], with hyper-parameter tuning specific for the MNI152 template. We then apply this non-linear transformation to our fMRI sequence in our T1w space. We apply the final transformation result only once to the EPI image, which prevents us from unnecessary fixed-precision multiplies, each of which will lead to a loss of information as we increase the number of multiplies due to the number of bits the numbers are represented as being insufficient.

### 3.1.3   Nuisance Correction

**General Linear Model**   Over the course of an fMRI scanning session, many sources of noise arise that must be corrected for in order to make quality data inferences. The scanner inadvertantly heats up (producing a high strength magnetic field for sessions lasting up to ten minutes produces an enormous amount of heat). As the scanner heats, the fine electronic equipment inadvertantly drifts in the signal it
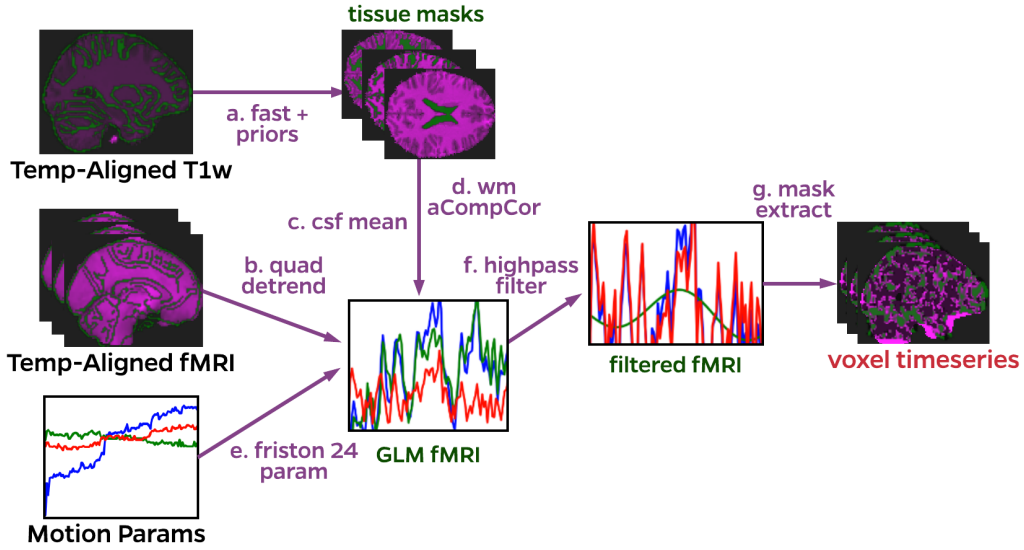
Figure 4: The nuisance correction workflow for the FNGS pipeline. We incorporate regressors for our Friston 24 motion parameters (from our original motion regressors), physiological regressors (from our csf mean signal and top 5 eroded white-matter principal components), as well as a quadratic drift regressors into a GLM. We remove the best-fit solution from our timeseries, as this signal can be thought of as being fit by regressor parameters effectively and therefore is likely not signal. The remaining signal is then high-pass filtered for fourier modes with a period over 100 seconds; our brain activity at resting state is on the order of around 16 seconds, and as long as we do not have tasks on the order of 100 seconds, we will not remove signal pertaining to our task condition. Finally, we extract our voxel timeseries from all volumes but the first 2 to avoid potential signal inconsistencies while the scanner is normalizing.

picks up (first demonstrated by [17] when they showed that a heated scanner detected "brain activity" in cadavers). This drift has been shown to be approximately quadratic [18].

While spatial motion correction removes the visual impact of head motion, spurious signal artifacts remain present. These artifacts can be characterized by the position of the brain in the scanner and the prior positions of the brain in the scanner, as first shown by Friston et al. [19]. This history relationship can be shown to be effectively captured by the current volume and the preceding volume, as well as their squares, so we introduce our friston 24 regressors where we have 4 regressors (1 current frame, 1 shifted frame, 1 squared-current frame, 1 square-shifted frame) for each of our 6 (x, y, z translation and rotation) motion regressors.

Finally, it has been shown by [20] that the fMRI signal is corrupted by physiological noise, from physiological functions such as blood flow or vessel dilation. The physiological confounds present in our functional data have been shown to be effectively captured by regressors for the mean lateral ventricle signal and the top 5 principal components from the white-matter signal [20], [21]. We estimate our csf and white-matter masks using the fast algorithm, [22] with priors obtained from the MNI152 parcellation [23]. This estimated mask is eroded by 2 voxels on all sides to avoid any potential signal distortion from the gray-matter signal, since gray-matter signal is expected to correlate with our stimuli, any signal bleeding into the white-matter voxels (since the gray-matter/white-matter boundary has a slight bleed-over region) that could get removed by our PCA would be disasterous on our downstream inferences.

To fit our regressors, we use a General Linear Model (GLM). For our $n$ voxels, the $t$ timestep BOLD signal, we can decompose $T_{raw} \in \mathbb{R}^{t \times n}$ as:

$$
T_{raw} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} c_{c,1} & c_{c,2} & \dots & c_{c,n} \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \\ \vdots \\ t \end{bmatrix} \begin{bmatrix} c_{l,1} & c_{l,2} & \dots & c_{l,n} \end{bmatrix} + \begin{bmatrix} 0 \\ 1 \\ \vdots \\ t^2 \end{bmatrix} \begin{bmatrix} c_{q,1} & c_{q,2} & \dots & c_{q,n} \end{bmatrix} + \epsilon
$$

$$
= \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 1 \\ \vdots & \vdots & \vdots \\ 1 & t & t^2 \end{bmatrix} \begin{bmatrix} c_{c,1} & c_{c,2} & \dots & c_{c,n} \\ c_{l,1} & c_{l,2} & \dots & c_{l,n} \\ c_{q,1} & c_{q,2} & \dots & c_{q,n} \end{bmatrix} + \epsilon
$$

$$
= RC + \epsilon
$$

where $\epsilon$ is the timeseries corrected without the quadratic $RC$. Since we know that brain dynamics behave non-quadratically, we know that minimizing the squared-error loss of $T$ with respect to $RC$ will

Figure 5: The connectivity matrix workflow for the FNGS pipeline.

provide a best-estimate of the coefficients $C$ of our quadratic regressors $R$, where our desired, quadratic-removed timeseries is just $\epsilon$ since we know that $\epsilon$ is the components of our raw signal that cannot be fit by a quadratic [18]. We can solve using the least-squares solution to regression problems:

$$C = (R^T R)^{-1} R^T T$$

Using our estimate $C$ for our regressors $R$, we can then find our quadratic-corrected timeseries to be:

$$\epsilon = T_{raw} - RC$$

**Low Frequency Drift Removal**   Using our quadratic-corrected timeseries, we can then remove low-frequency drift that may still be present in our data. We know that any physiological response due to a stimulus will have a period of around 16 seconds (or a frequency of 0.2 Hz) and not exceeding the period of any stimuli present, as it has been shown in [24]. Using this information, we can select to remove sinusoidal components with frequencies far exceeding those of our predicted responses. Conservatively, we set a threshold of 0.01 Hz for highpass-filtering out low-frequency noise (this should not remove task-dependent signal as long as our task has a period less than about 100 seconds).

**T1 Effect Removal**   During the fMRI session, the first few volumes may appear to have brighter intensities as the T1 effects are not fully saturated. External attempts that remove the T1 effect include component-correction and global mean normalization of each slice, both of which have been shown to potentially remove brain signal [25]. To account for this, we simply discard the first 2 volumes from our fMRI sequence, which tends to account for the majority of contrast-dependent issues.

### 3.1.4   Connectivity Estimation

**Voxel Timeseries**   In fMRI processing, we often want to think about each voxel in the brain itself as a function of the time spent in the scanner. This allows us the ability to perform further analyses leveraging the temporal and spatial characteristics of each voxel. To accomplish this, we take a mask of known brain voxels associated with the registration template, and we extract the voxels of known brain tissue.

**ROI Timeseries** Similarly, it is often of neurological significance to think of the brain as segmented into regions of neurons performing a similar task based on their spatial layout. To accomplish this, neuroscientists parcellate the registered brains into parcellation atlases, whereby they characterize individual voxels in the registration template into individual functional or structural groups. For each region in the parcellation atlas, we subset all of the voxels associated with the given region of interest (ROI) and average them spatially. This gives us a significantly downsampled representation of the brain, often yielding in excess of 100 fold compression and also reducing voxel-by-voxel noise that may still be present. The downsampled timeseries can then be analyzed using traditional statistical methods with greater robustness and computational efficiency [26].

Finally, we can use our ROI timeseries to estimate a connectome. For our connectome, we want each edge to represent a relationship between two regions of the brain. Intuitively, the functional timeseries gives a functional relationship, whereby the simplest comparison is in terms of the firing simultaneity of two regions. For this, we estimate the correlation between each pair of regions from our ROI timeseries for each region in the brain. This yields us a simplified network representation of our original brain.

## 3.2 Discriminability

Throughout our investigation, we assume that our data is structured as:

$$T_{n,t} = g_{\psi,t}(f_{\varphi,t}(v_n))$$

Thus, we take an explicit observation $T_{n,t}$ of some latent signal $v_n$ for subject $n$. However, this signal is distorted, first by measurement distortion $f_{\varphi,t}$ and second by the processing options we have chosen $g_{\psi,t}$. Here, $\varphi$ represents error introduced by scanner parameters, and $\psi$ represents error introduced by our processing options themselves at a particular measurement session $t$. As the measurement and processing-dependent distortions are random and unknown, we intuitively want a robust pipeline to be one in which $T_{n,t}$ is as close as possible to $T_{n,t'}$, or the measurement of subject $n$ taken at a time $t'$. To accomplish this, we measure the discriminability $r(\psi, \varphi)$ [6], which provides a statistic that is near 1 when the observations $T_{n,t}$ and $T_{n,t'}$ are similar, and near 0.5 when $T_{n,t}$ and $T_{n',t'}$ for a different subject $n'$ are indistinguishable. Then we seek the pipeline that maximizes:

$$\psi_{best}, \varphi_{best} = \underset{\psi, \varphi}{\operatorname{argmax}} \, r(\psi, \varphi)$$

During a resting-state MRI session, researchers do not control the stimuli leading to brain responses in time. Theoretically, the idea is to measure the resting-state fluctuations unique to a particular person,

which may or may not be the exact same from one session to the next. Then comparing timeseries to timeseries does not make experimental sense; we have no idea what stimuli are impacting the observable dynamics of the brain at a given time, so we cannot assume that they are the same from one session to the next. To overcome this obstacle, we instead assume that, over the course of a scanning session, we get enough variation of a person's brain activity to be able to make inferences about which regions in the brain function together. We tune FNGS with respect to the maximal discriminability with the scan id as the label. This provides an upper-bound on downstream inference errors, by ensuring maximum similarity between observations of the same subject that are thought to be maximally similar [6].

To estimate this functional connectivity, we try several strategies:

**Correlational Approach** The simplest and most intuitive functional connectome is the connectivity matrix as we have previously described and provide natively as the connectome estimation in the FNGS pipeline. Here, we consider the correlation over time between each pair of edges in our ROI timeseries [27]. With this definition of a spatially-motivated functional connectome, we assume that over time, regions that are more connected will tend to fire simultaneously in a similar way from session to session.

**Fourier Approach** Moreover, it might also be valid to think of the timeseries in terms of its fourier properties. While temporal properties themselves may not be consistent from session to session, instead it may be the case that the frequency with which individual ROI timeseries opearte is more characteristic of a given person's brain activity [28]. Following Chen et. al, we first consider the fourier transform for each ROI, and highpass filter any activity below the 0.01 Hz range, as this activity can be thought of as low-frequency drift. We normalize the remaining fourier components to form a probability distribution. Next, we consider the Kullback-Leibler divergence between each pair of regions for the amplitude or power spectrum to obtain a spatial connectome for each subject that instead considers a notion of the frequency-domain similarities between regions.

**Ranking** In addition, for both the correlational and fourier approaches above, we consider the ranked approximation of each connectome. To compute a ranked connectome, we rank each edge in the connectivity matrix from lowest to highest value with a scalar ranking, breaking ties with the mean rank. Essentially our hypothesis is that this allows relative robustness, is affine-transformation invariant (ie, a global scaling of the edge weights will not be reflected after ranking), and allows us to recover monotonic nonlinear fluctuations that may be present better than a linear-transformation such as z-scoring of the correlations [6].

The distance matrix is required for computation of the discriminability score from a connectome.

**Frobenius Norm of Difference** For our correlational estimates, we consider the distance matrix $D$ to be:

$$D(A, B) = ||A - B||_F$$

or the frobenius norm of the difference each pair of connectomes $A, B$ obtained in our study, where $A, B$ are correlation-derived connectomes.

**Hellinger Distance** For our frequency-derived connectomes, we consider the distance matrix $H$ to be:

$$H(C, D) = \frac{1}{\sqrt{2}} ||\sqrt{C} - \sqrt{D}||_2$$

or the hellinger distance between each pair of connectomes $C, D$ obtained in our study, where $C, D$ are fourier-derived connectomes.

## 3.3 Software

### 3.3.1 Docker Container

Through FNGS, we have developed a reliable functional connectome estimation solution. To maximize the deployability of our software,we have developed a publicly-accessible docker container preloaded with all of our software. Using docker ensures compatibility with any modern (Windows, Linux, OSx) computing platform, as all dependencies are folded into the container directly.

### 3.3.2 Deployment

Leveraging our docker container, we are able to offer multiple levels of access for the FNGS pipeline. The pipeline can directly analyze single-scans using any input format the researcher chooses with a local entrypoint in the docker container. For analysis on large numbers of scans at once, we require users to format their data according to the BIDs spec [29], a popular structure for formatting fMRI, diffusion MRI, and various other modalities of neuroimaging data. This option supports parallel performance by setting options for the number of threads available. Typically, this will be ablout $min\left(ncpus, \frac{max(RAM)}{6}\right)$ to be conservative (the pipeline generally sues about 3-4 gigs of RAM max, so this gives us some headroom). Users can also optionally call the FNGS pipeline on their local data using another entrypoint in the docker container. If users do not have a high performance computer or a cluster, we leverage Amazon Web Services [30] to allow users to batch-analyze large amounts of data for relatively low cost. We provide a command-line entrypoint and automation scripts that allows users to scale their analyses in parallel on the cloud with virtually no setup required other than appropriate configuration of their Amazon Web
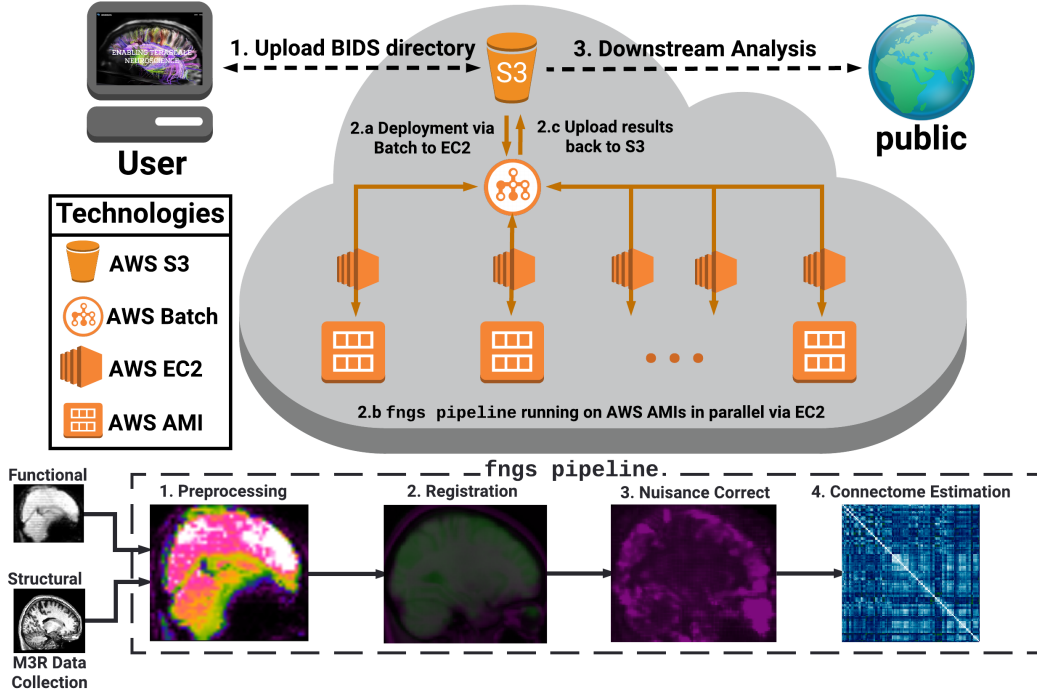
Figure 6: The deployment workflow for the FNGS pipeline. Users provide a directory according to the BIDs spec to the FNGS cloud controller remotely via the FNGS web-service or locally via the docker container at (1). The data is uploaded directly to AWS S3 cloud drives. The controller then initiates the Batch deployment procedure at (2.a), which interfaces between S3 and EC2 cloud computers to provide the MRI scans to EC2 instances pre-loaded with the FNGS pipeline for analysis. After the scans are finished being analyzed on the EC2 instances at (2.b), the results are then re-uploaded back to the S3 cloud drive at (2.c). The user can then navigate to the S3 cloud drive for downstream analyses.

Services account. The procedure can be seen in Figure (6) and is detailed in our provided tutorials below.

### 3.3.3 Resources

Tutorials for usage and demos, and links to the code, and our live version of the website can be found at https://github.com/neurodatadesign/fngs. All of our software is open source and coded in python. An exhaustive high-level overview of the pipeline and all derivatives can be found at:

https://neurodatadesign.github.io/fngs/about_fngs/Schematic.html. Note that the image is interactive.

We develop an R package with many utilities for easily computing and comparing discriminability here: https://github.com/ebridge2/Discriminability.

# References

[1] S Ogawa, T M Lee, A R Kay, and D W Tank. Brain magnetic resonance imaging with contrast dependent on blood oxygenation. *Proceedings of the National Academy of Sciences of the United States of America*, 87(24):9868–72, dec 1990.

[2] Richard B. Buxton, Kâmil Uludağ, David J. Dubowitz, and Thomas T. Liu. Modeling the hemodynamic response to brain activation. *NeuroImage*, 23:S220–S233, jan 2004.

[3] K. J. Friston, C. D. Frith, P. F. Liddle, and R. S. Frackowiak. Functional connectivity: the principal-component analysis of large (PET) data sets. *Journal of cerebral blood flow and metabolism*, 13(1):5–14, January 1993.

[4] Urs Braun, Sarah F Muldoon, Danielle S Bassett, Urs Braun, Sarah F Muldoon, and Danielle S Bassett. On Human Brain Networks in Health and Disease. In *eLS*, pages 1–9. John Wiley & Sons, Ltd, Chichester, UK, feb 2015.

[5] Stephen Strother, Stephen La Conte, Lars Kai Hansen, Jon Anderson, Jin Zhang, Sujit Pulapura, and David Rottenberg. Optimizing the fMRI data-processing pipeline using prediction and reproducibility performance metrics: I. A preliminary group analysis. *NeuroImage*, 23:S196–S207, 2004.

[6] Shangsi Wang, Zhi Yang, Xi-Nian Zuo, Michael Milham, Cameron Craddock, Gregory Kiar, William Gray Roncal, Eric Bridgeford, Carey E Priebe, and Joshua T Vogelstein. Optimal Decisions for Discovery Science via Maximizing Discriminability: Applications in Neuroimaging. *request for preprint*.

[7] Nathan W Churchill, Anita Oder, Hervé Abdi, Fred Tam, Wayne Lee, Christopher Thomas, Jon E Ween, Simon J Graham, and Stephen C Strother. Optimizing Preprocessing and Analysis Pipelines for Single-Subject FMRI. I. Standard Temporal Motion and Physiological Noise Correction Methods. *Human Brain Mapping*, 33:609–627, 2012.

[8] K J Gorgolewski, N Mendes, D Wilfling, E Wladimirow, C.J. Gauthier, T Bonnen, R Trampel, P L Bazin, and D S Margulies. Measuring variability of human brain activity at rest - a high resolution 7-Tesla test-retest dataset.

[9] Mark W. Woolrich, Saad Jbabdi, Brian Patenaude, Michael Chappell, Salima Makni, Timothy Behrens, Christian Beckmann, Mark Jenkinson, and Stephen M. Smith. Bayesian analysis of neuroimaging data in FSL. *NeuroImage*, 45(1):S173–S186, 2009.

[10] Koene R A Van Dijk, Mert R Sabuncu, and Randy L Buckner. The influence of head motion on intrinsic functional connectivity MRI. *NeuroImage*, 59(1):431–8, jan 2012.

[11] Mark Jenkinson, Peter Bannister, Michael Brady, and Stephen Smith. Improved optimization for the robust and accurate linear registration and motion correction of brain images. *NeuroImage*, 17(2):825–41, oct 2002.

[12] Donald G. Mitchell and Mark Cohen. *MRI principles*. Saunders, 2004.

[13] M Jenkinson and S Smith. A global optimisation method for robust affine registration of brain images. *Medical image analysis*, 5(2):143–56, jun 2001.

[14] Douglas N. Greve and Bruce Fischl. Accurate and robust brain image alignment using boundary-based registration. *NeuroImage*, 48(1):63–72, oct 2009.

[15] Günther Grabner, Andrew L. Janke, Marc M. Budge, David Smith, Jens Pruessner, and D. Louis Collins. Symmetric Atlasing and Model Based Segmentation: An Application to the Hippocampus in Older Adults. pages 58–66. Springer, Berlin, Heidelberg, 2006.

[16] Jesper L R Andersson, Mark Jenkinson, Stephen Smith, and Jesper Andersson. Non-linear registration aka Spatial normalisation. 2007.

[17] Anne M Smith, Bobbi K Lewis, Urs E Ruttimann, Frank Q Ye, Teresa M Sinnwell, Yihong Yang, Jeff H Duyn, and Joseph A Frank. Investigation of Low Frequency Drift in fMRI Signal. 1999.

[18] Jody Tanabe, David Miller, Jason Tregellas, Robert Freedman, and Francois G Meyer. Comparison of Detrending Methods for Optimal fMRI Preprocessing.

[19] Karl J. Friston, Steven Williams, Robert Howard, Richard S. J. Frackowiak, and Robert Turner. Movement-related effects in fmri time-series. *Magnetic Resonance in Medicine*, 35(3):346–355, 1996.

[20] Y Behzadi, K Restom, J Liau, and T T Liu. A component based noise correction method (CompCor) for BOLD and perfusion based fMRI. *Neuroimage*.

[21] Rastko Ciric, Daniel H. Wolf, Jonathan D. Power, David R. Roalf, Graham L. Baum, Kosha Ruparel, Russell T. Shinohara, Mark A. Elliott, Simon B. Eickhoff, Christos Davatzikos, Ruben C. Gur, Raquel E. Gur, Danielle S. Bassett, and Theodore D. Satterthwaite. Benchmarking of participant-level confound regression strategies for the control of motion artifact in studies of functional connectivity. *NeuroImage*, 154:174 – 187, 2017. Cleaning up the fMRI time series: Mitigating noise with advanced acquisition and correction strategies.

[22] Y. Zhang, M. Brady, and S. Smith. Segmentation of brain MR images through a hidden Markov random field model and the expectation-maximization algorithm. *IEEE Transactions on Medical Imaging*, 20(1):45–57, jan 2001.

[23] G Grabner, A L Janke, MM Budge, D Smith, J Pruessner, and D L Collins. Symmetric atlasing and model based segmentation: an application to the hippocampus in older adults. *Med Image Comput Comput Assist Interv Int Conf Med Image Comput Comput Assist Interv*, 9, 2006.

[24] Martin A Lindquist, Ji Meng Loh, Lauren Y Atlas, and Tor D Wager. Modeling the hemodynamic response function in fMRI: efficiency, bias and mis-modeling. *NeuroImage*, 45(1 Suppl):S187–98, mar 2009.

[25] Molly G Bright, Christopher R Tench, and Kevin Murphy. Potential pitfalls when denoising resting state fMRI data using nuisance regression. 2016.

[26] Russell A. Poldrack. Region of interest analysis for fMRI. *Social Cognitive and Affective Neuroscience*, 2(1):67–70, mar 2007.

[27] Baxter P Rogers, Victoria L Morgan, Allen T Newton, and John C Gore. Assessing functional connectivity in the human brain by fMRI. *Magnetic resonance imaging*, 25(10):1347–57, dec 2007.

[28] Cencheng Shen, Carey E Priebe, Mauro Maggioni, and Joshua T Vogelstein. Discovering the Geometry of Relationships Across Disparate Data Modalities. *submitted*.

[29] Krzysztof J. Gorgolewski, Tibor Auer, Vince D. Calhoun, R. Cameron Craddock, Samir Das, Eugene P. Duff, Guillaume Flandin, Satrajit S. Ghosh, Tristan Glatard, Yaroslav O. Halchenko, Daniel A. Handwerker, Michael Hanke, David Keator, Xiangrui Li, Zachary Michael, Camille Maumet, B. Nolan Nichols, Thomas E. Nichols, John Pellman, Jean-Baptiste Poline, Ariel Rokem, Gunnar Schaefer, Vanessa Sochat, William Triplett, Jessica A. Turner, Gaël Varoquaux, and Russell A. Poldrack. The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Scientific Data*, 3:160044, jun 2016.

[30] Amazon Inc. *Amazon Elastic Compute Cloud (Amazon EC2)*. Amazon Inc., 2008.