



Common Data Elements, Scalable Data Management Infrastructure, and Analytics Workflows for Large-Scale Neuroimaging Studies

Rayus Kuplicki^{1*}, James Touthang¹, Obada Al Zoubi¹, Ahmad Mayeli¹, Masaya Misaki¹, NeuroMAP-Investigators^{1,2}, Robin L. Aupperle^{1,2}, T. Kent Teague^{3,4,5}, Brett A. McKinney^{6,7}, Martin P. Paulus¹ and Jerzy Bodurka^{1,8*}

OPEN ACCESS

Edited by:

Martin Walter,
University of Tübingen, Germany

Reviewed by:

Martin Dyrba,
German Center for
Neurodegeneratives, Helmholtz
Association of German Research
Centers (HZ), Germany
Rajat Mani Thomas,
Academic Medical
Center, Netherlands

*Correspondence:

Rayus Kuplicki
rkuplicki@laureateinstitute.org
Jerzy Bodurka
jbodurka@laureateinstitute.org

Specialty section:

This article was submitted to
Neuroimaging and Stimulation,
a section of the journal
Frontiers in Psychiatry

Received: 19 March 2021

Accepted: 19 May 2021

Published: 17 June 2021

Citation:

Kuplicki R, Touthang J, Al Zoubi O,
Mayeli A, Misaki M,
NeuroMAP-Investigators,
Aupperle RL, Teague TK,
McKinney BA, Paulus MP and
Bodurka J (2021) Common Data
Elements, Scalable Data Management
Infrastructure, and Analytics
Workflows for Large-Scale
Neuroimaging Studies.
Front. Psychiatry 12:682495.
doi: 10.3389/fpsy.2021.682495

¹ Laureate Institute for Brain Research, Tulsa, OK, United States, ² Department of Community Medicine, Oxley College of Health Sciences, University of Tulsa, Tulsa, OK, United States, ³ Department of Surgery, University of Oklahoma School of Community Medicine, Tulsa, OK, United States, ⁴ Department of Psychiatry, University of Oklahoma School of Community Medicine, Tulsa, OK, United States, ⁵ Department of Biochemistry and Microbiology, Oklahoma State University Center for Health Sciences, Tulsa, OK, United States, ⁶ Department of Mathematics, University of Tulsa, Tulsa, OK, United States, ⁷ Tandy School of Computer Science, University of Tulsa, Tulsa, OK, United States, ⁸ Stephenson School of Biomedical Engineering, University of Oklahoma, Norman, OK, United States

Neuroscience studies require considerable bioinformatic support and expertise. Numerous high-dimensional and multimodal datasets must be preprocessed and integrated to create robust and reproducible analysis pipelines. We describe a common data elements and scalable data management infrastructure that allows multiple analytics workflows to facilitate preprocessing, analysis and sharing of large-scale multi-level data. The process uses the Brain Imaging Data Structure (BIDS) format and supports MRI, fMRI, EEG, clinical, and laboratory data. The infrastructure provides support for other datasets such as Fitbit and flexibility for developers to customize the integration of new types of data. Exemplar results from 200+ participants and 11 different pipelines demonstrate the utility of the infrastructure.

Keywords: human brain, neuroimaging, multi-level assessment, large-scale studies, common data element, data processing pipelines, scalable analytics, bids format

INTRODUCTION

Neuroimaging studies such as ABCD, ADNI, Human Connectome, and Tulsa 1,000 studies are significant contributors to the rapid growth of big data (1–4). In addition to the usual high-dimensional data that accompany clinical studies (e.g., genetic, cellular, and clinical assessments), neuroscience studies include multimodal data for the brain [e.g., MRI, Perfusion MRI [pMRI], diffusion MRI [dMRI], functional MRI [fMRI] and Electroencephalography [EEG]]. The use of various data acquisition modalities and differences in studies' experimental designs make it challenging to provide a common data architecture that would offer easy access, scalability, management and sharing, including the ability to build analytic workflows and to run large scale analyses with increasingly large numbers of subjects. Here, we propose possible solutions to these challenges and described our specific working implementation.

As a part of the Neuroscience-Based Mental Health Assessment and Prediction (NeuroMAP) Center of Biomedical Research Excellence (CoBRE) award from National Institute of General Medical Sciences (NIGMS/NIH), the NeuroMAP Research Core provides research infrastructure

to conduct advanced neuroscience research and is also responsible for providing active data management and analysis support, which includes standardization of all acquired data. Data collected for NeuroMAP consist of a core baseline assessment as well as subsequent individual projects sharing various common data elements. Briefly, the research core protocol contains neuroimaging (two sessions—one functional with concurrent EEG lasting 2 h and one structural lasting 1 h), behavioral, self-report, biomarker, and actigraphy data acquired from large cohorts of participants who are then enrolled in the various other projects. Full details can be found in the supplement, especially see **Supplementary Figure 1** for an overview of the core. Ongoing human recruitment into the core protocol is roughly 100 participants per year in phase I (5 years, with a possible extension to 10 years), so that this cohort is anticipated to reach 400+ participants. Currently at year 3, 310 participants have been enrolled, with 291 completing all core assessments (see **Supplementary Table 1** for general sample characteristics). A large and growing cohort size combined with several acquisition modalities amounts to a large and increasing set of heterogeneous and complex data.

Large-scale data collection pipelines are complex to establish while maintaining standardized experimental protocols on both the data-acquisition hardware level and on the clinical data management level. Follow-up analyses also require further standardization, which is often implemented in *ad hoc* software systems at different institutions and may even vary between labs within an institution. Home-grown solutions can work adequately, and over the past decade we have collected neuroimaging data from thousands of individuals using our own internal solutions. However, in recent years, progress has been made in the scientific community toward consensus solutions to improve data management and mechanisms for data sharing (5).

There are a number of substantial costs when using custom data management solutions, not the least of which is developing the data processing standards, which can be difficult for researchers without informatics training. Practically speaking, the naming conventions and processing steps used in a study are often neither well-documented nor reproducible. In the best case, idiosyncratic naming conventions and directory structures simply add overhead when sharing datasets and analysis code that was developed for specific file structures. For example, a researcher unfamiliar with a particular dataset would need to learn about its conventions along with the details of the study. Sometimes, the first thing researchers do when working with a new dataset is reformat it to match a form they are familiar with, which is extra effort that could be avoided if standard formats were used. Similarly, reusing analysis code (e.g., scripts and software) often requires either extensive reworking to be compatible with a new dataset, or reformatting the target data to be compatible with the existing code.

One possible solution is the development of a complex data management system used to store, access, and even analyze neuroimaging and associated data. There have been several projects to produce such extensive systems over the past 15 years (6–12); however, they can come with significant overhead in installation, maintenance, and user training. In fact, our institute

spent considerable time and resources attempting to implement one of these systems, a project which we ultimately abandoned due to excessive cost and technical difficulties.

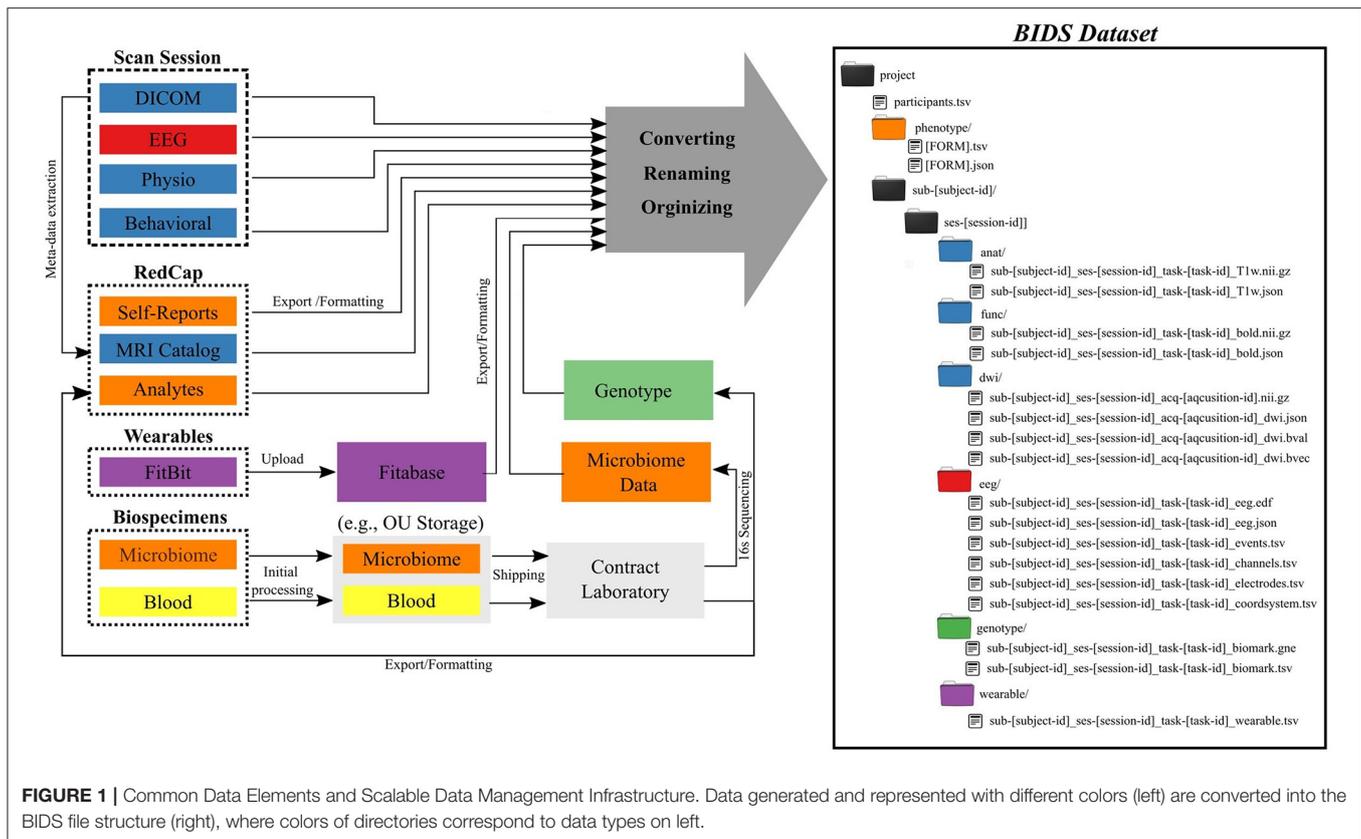
One of the main challenges is the need for a commonly accepted data structure format that would provide a consistent and standardized way to organize multi-level neuroimaging data. The Brain Imaging Data Structure (BIDS) (5) was introduced in 2016 and promises to alleviate some of the difficulties in organizing, documenting and sharing data and code while maintaining a simple, intuitive structure that is easy to understand and work with. With metadata stored directly on disk, either in the form of file names and locations or associated JSON sidecars, BIDS avoids requiring overly complex management software or databases. The BIDS format is remarkably similar to our internally developed neuroimaging data organization solution and we decided to transition to BIDS for the NeuroMAP studies, common data elements and all new projects going forward. Wide acceptance of BIDS provides standardization across other datasets and facilitates sharing with the scientific community.

The goal of this work is to provide a detailed description of our computing and data management infrastructure, which will contribute to common data structures and serve as a model for other large studies and institutions. We also show results from a few exemplar datasets/analyses as a proof of concept. The rest of the manuscript is organized as follows: Methods describes our architecture and workflows. Data Management Infrastructure Design provides an overview of our overall infrastructure, with BIDS Conversion containing details related to organization of five different modalities and Analytic Workflows detailing the analysis pipeline structure. Results shows some exemplar group results from task and resting fMRI as well as EEG. Discussion summarizes our current state and future directions.

METHODS

Data Management Infrastructure Design

The Common Data Elements and Scalable Data Managing Infrastructure can integrate neuroimage data with various other data types (**Figure 1**). The CDE data are in general composed from multimodal MRI, fMRI, EEG, physiological recordings, behavioral measures, self-reports measures, actigraphy from wearable devices, and biospecimen samples (e.g., blood and microbiome). All self-reported and clinical interview data are collected directly into an internally-hosted instance of the Research Electronic Data Capture (REDCap) (13) database. REDCap is a secure (e.g., compliant with 21 CFR Part 11, FISMA, HIPAA, and GDPR), web-based data collection system that is currently used by over 3,200 institutions in 128 countries (14). We adopted REDCap in 2014 to replace a combination of paper-charts and a home-grown databasing solution. Full details of the NeuroMAP core multilevel data collection are included in **Supplementary Materials**. All original data sources (left side of **Figure 1**) are processed and stored in order to produce a BIDS-compliant dataset (right side of **Figure 1**). The middle part of the figure shows intermediate steps and storage, while the right shows the final BIDS dataset. BIDS conversion of each element is



described in detail in Section BIDS Conversion. Colors are used to show which raw data and samples correspond to particular elements of the BIDS dataset in its final form.

BIDS Conversion

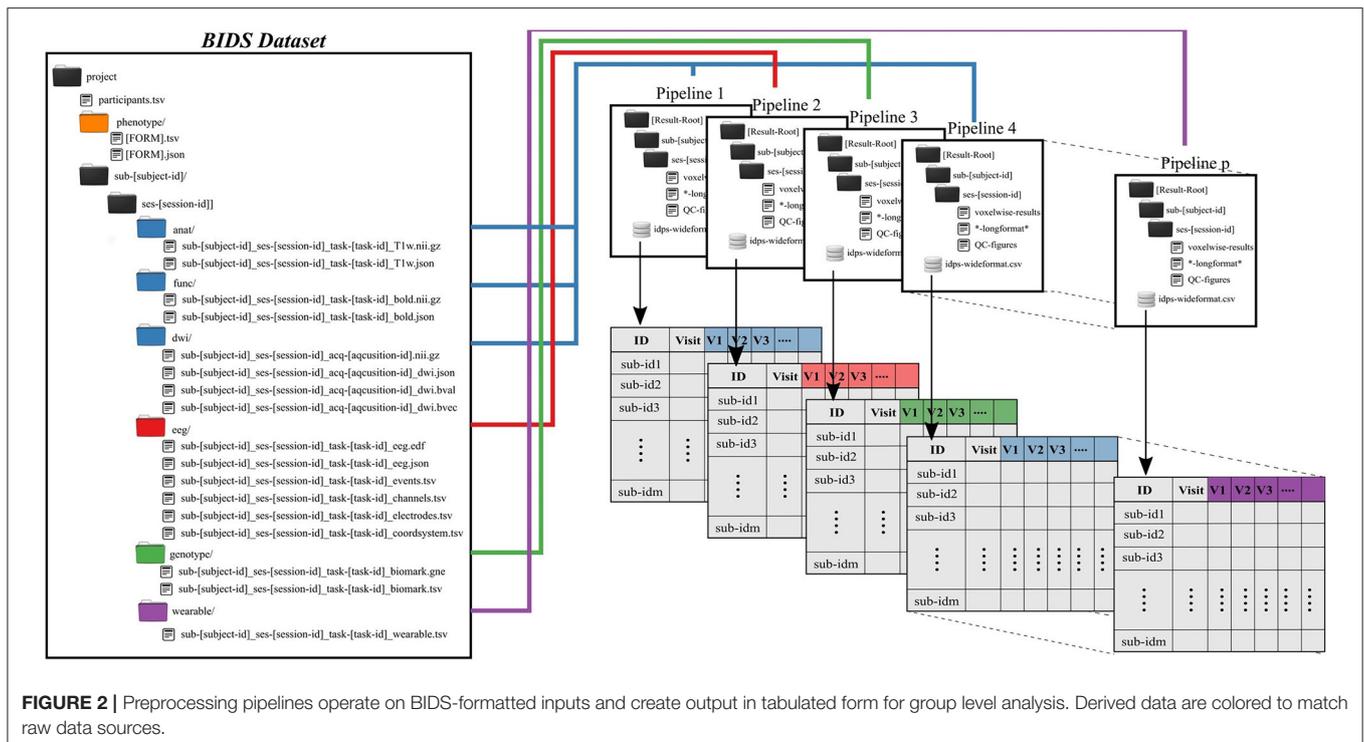
Self-Report/REDCap

Self-report and clinical measures (described in full in section S-1.1.1) stored in REDCap are exported into a BIDS-compliant format using the PyCap library built on top of the REDCap API. The inputs/outputs of this process appear in orange in **Figure 1**. In brief, an API key links a user and access rights to a single project. Data returned from REDCap include a table of subject data for the project as well as metadata about the project and data collection instruments. The data are converted to tsv format and stored in the phenotype folder following BIDS specification. Similarly, the metadata describing the data collection instruments are stored in JSON formatted data dictionaries. The result is a JSON/tsv pair for each REDCap form. This script can be setup for other redcap projects and is available on GitHub (<https://github.com/laureate-institute-for-brain-research/redcap-to-bids>).

Neuroimaging and Associated Physiological Data

This section describes the collection and organization of neuroimaging and associated data across all studies. For specific details regarding the NeuroMAP Core, see section S1.2 of the supplement. Neuroimaging data are produced in two formats.

Source DICOM images are reconstructed and generated by the scanner and permanently stored in a read-only central location. The default organization from GE DICOM file structure has each scan stored three-folders deep (e.g., pXXX/eYYY/sZZZ, where p, e, and s refer to patient, exam, and series). For each completed scan and patient exam, these DICOM images are automatically extracted, transferred to scanner-dedicated local storage and reorganized by custom developed real-time MRI scanner data management software. To reduce the storage burden associated with hundreds of thousands of individual files, DICOM folders are packaged in.tar.gz format at the exam directory level. This reduces the number of individual files stored by a factor of 10^5 , and also saves significant storage space when individual files are smaller than the storage block size. Each DICOM image contains standard metadata indicating the subject ID, date, study, scan, and various imaging parameters: everything necessary to associate a scan with its final BIDS-compliant name and location. However, parsing through the DICOM folders and extracting metadata is an expensive operation, even before considering the compressed format. We solved this problem by creating a REDCap project called the MRI Catalog, which contains all relevant DICOM metadata. New DICOM images from MRI scans are processed and metadata describing them are imported into REDCap nightly. Our real-time MRI software also produces a unique exam folder (on scanner-attached and dedicated real-time processing Linux workstations), which contains AFNI formatted imaging data that are uploaded and created in real



time from a given session, along with any associated concurrent physiological recordings (pulse oximeter, respiratory belt, pre-processed EEG), electronic documentation for each scan with imaging parameters, DICOM file count and location on the local storage after extraction from the MRI scanner host computer and image database.

Raw EEG data (without any preprocessing) acquired concurrent with fMRI are initially stored locally on a dedicated EEG recording computer and then synchronized and transferred to network storage nightly. Similarly, behavioral responses collected during scanning tasks are initially stored on a stimulus laptop and then moved to network storage immediately upon session completion. The decision to store data locally first, then move it to network storage was based on reliability and latency considerations, so that networking issues do not affect data collection.

Neuroimaging and associated physiological data are organized and converted to BIDS format by a nightly batch process. This process handles the neuroimaging and behavioral data separately. In the first step, an export of all current MRI Catalog data necessary for organization is extracted from REDCap. The organization process parses through these data looking for project and scan IDs matching lists for a particular project. Newly acquired matching scans are converted to nii.gz format and sent to the appropriate BIDS folder with an associated JSON sidecar. Importantly, the DICOM metadata also contains a pointer to the appropriate exam folder and series number, which is used to extract the associated physiological data. Technical issues often make data collection imperfect, e.g., scans may be aborted/restarted due to participant discomfort or imaging

artifacts. Therefore, quality checks take place to help maintain data fidelity. The two most relevant checks include subject and duration matching. REDCap contains a list of subjects who have been consented for each study, so any subject ID in the MRI Catalog that does not match a consented subject for the study in question is not included. This happens, for example, with technical scans, which should not appear in the final dataset. The case where scans are repeated, producing multiple scans of the same type is handled by matching on expected duration. Any scan that does not have the expected duration is discarded, since shortened duration indicates an incomplete scan.

The second part of the organization process handles new behavioral data found on network storage. These data are stored in a folder unique to the study and completion date/time of the session. Each behavioral folder should contain data from one subject at one visit, and any folders that contain multiple subjects or visits generate an error and are skipped until they are manually corrected. Reformatting raw behavioral data involves converting from csv to tsv, creation of a new header, and then placement in the final BIDS data structure. Raw EEG data are named according to subject ID, and quality control involves matching on subject ID, date/time, and duration, similar to what is done for imaging data.

Behavioral Data

Data management for behavioral sessions completed outside the scanner mirrors that for the behavioral data from scanning sessions, where raw files are initially stored locally, moved to network storage at the end of the session, and then parsed/organized nightly. The behavioral session also

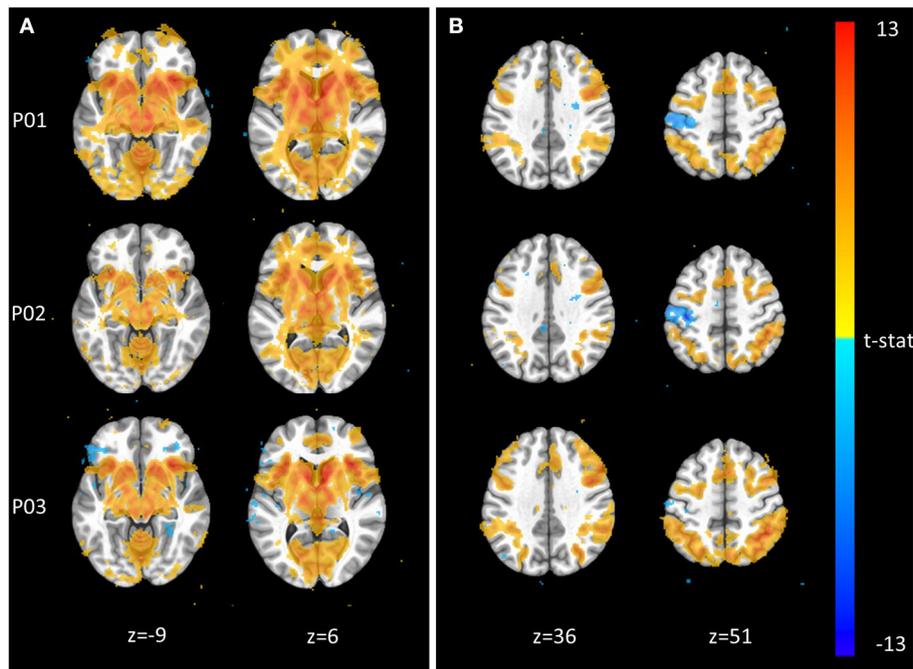


FIGURE 3 | Exemplar voxel-wise task activation maps produced by three different pipelines. **(A)** Monetary Incentive Delay P5–P0 contrast from $n = 93$ participants at $p < 0.001$. **(B)** Stop Signal Stop–NoStop contrast from $n = 49$ subjects at $p < 0.001$.

includes physiological data acquired using Acknowledge software (BIOPAC Systems, Inc.). These data are initially stored as a single continuous file in.acq format covering the entire session. Bioread (<https://github.com/uwmadison-chm/bioread>) is used to convert to plain text format, which is then sliced into and saved as individual tsv.gz files for each task and run. Synchronization is done using the parallel port, with a unique code indicating the start and end of each task. The appropriate header values are also extracted and stored in a JSON sidecar to be BIDS compliant.

Biospecimens

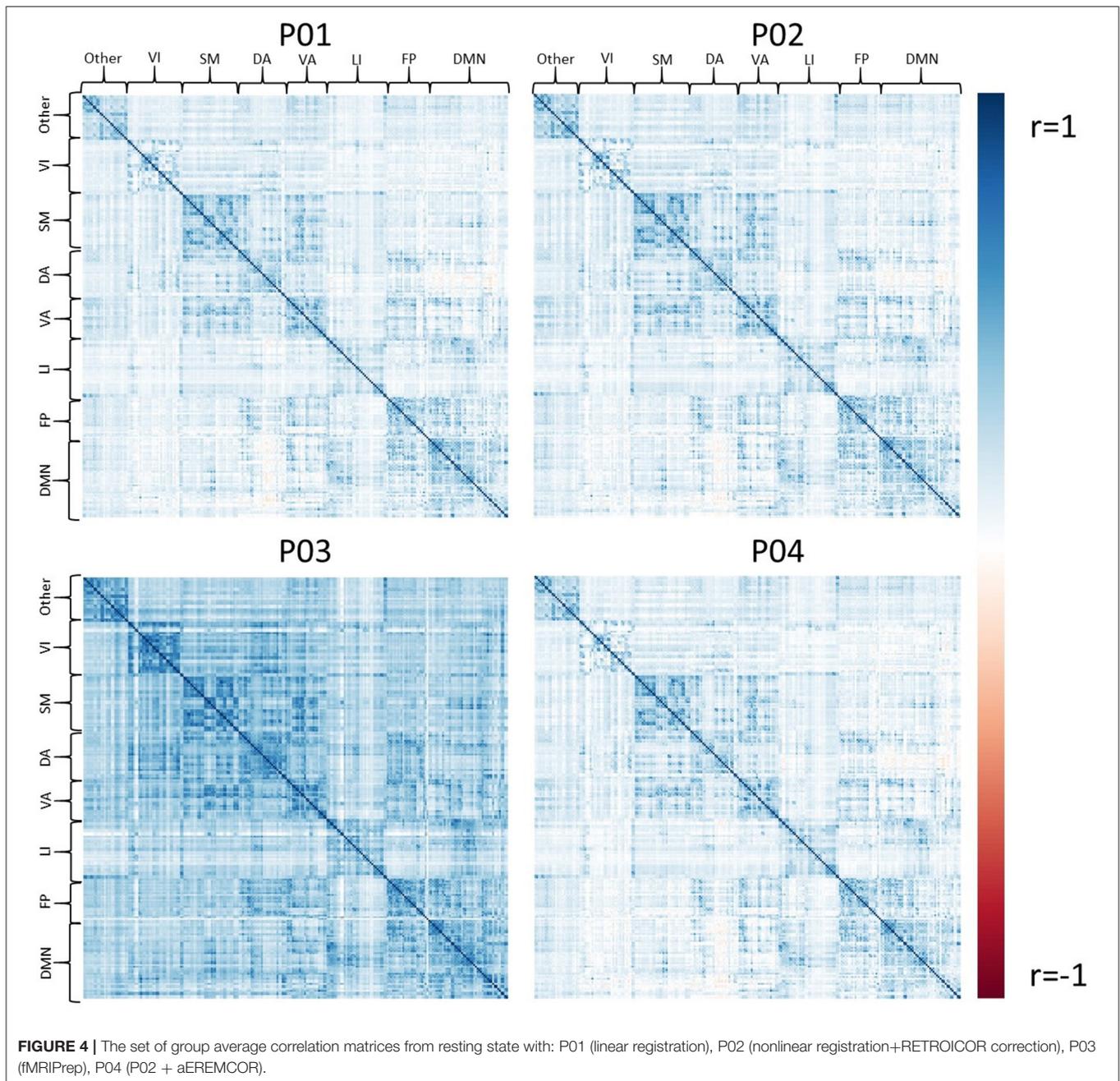
A detailed description of initial processing and storage of biospecimens is in the supplement (S 1.4). Final processing of the collected samples may be carried out by a contract laboratory or done in-house and produces datasets of varying size. Blood samples are used to quantify a limited number of analytes (e.g., <50) describing inflammatory and metabolic states. These data are parsed and imported into REDCap for permanent storage, and then later exported into BIDS format in the same way as self-report scales. Blood samples are also sent for genotyping, which produces 650,000 or more values per participant. These data are not suitable for storage in REDCap, so they are stored in a separate repository where the location and genetic descriptors are identified in the BIDS data description. Microbiome samples produce similarly large datasets through 16S sequencing or other technologies, which again are identified in the data description to be BIDS compliant and do not have permanent storage within REDCap.

Actigraphy/FitBit

FitBit data are initially stored in a third-party database (Fitabase <https://www.fitabase.com/>, accessed 2/18/2021), which handles most of the overhead related to FitBit account creation/management and aggregation of many participants' data. Data exported from Fitabase may be divided into daily summaries and momentary assessments. Due to account management details, daily summary data often include time periods outside of the assessment windows for each subject. Start and end dates, entered into REDCap by the researcher deploying the FitBit, are used to trim the summary data down to the appropriate timeframe. These daily summaries are stored in a single table under the phenotype folder and include overall activity levels, sleep duration and quality. Momentary assessment data including minute-wise heart rate estimates are stored in each subject's wearable folder and are in many ways similar to behavioral outputs. Fitabase provides FitBit data in four different time intervals: 30 s, 1 min, 1 h, and 24 h. Thirty-second interval data only includes sleep stages. Minute interval data include calories burned, activity intensity, metabolic equivalent of tasks (METs), current sleep stage, heart rate, and number of steps. One-hour interval data include calories burned, activity intensity, and number of steps. Twenty-four-hour interval data include activity summaries, calories burned, number of steps, and sleep.

Analytic Workflows

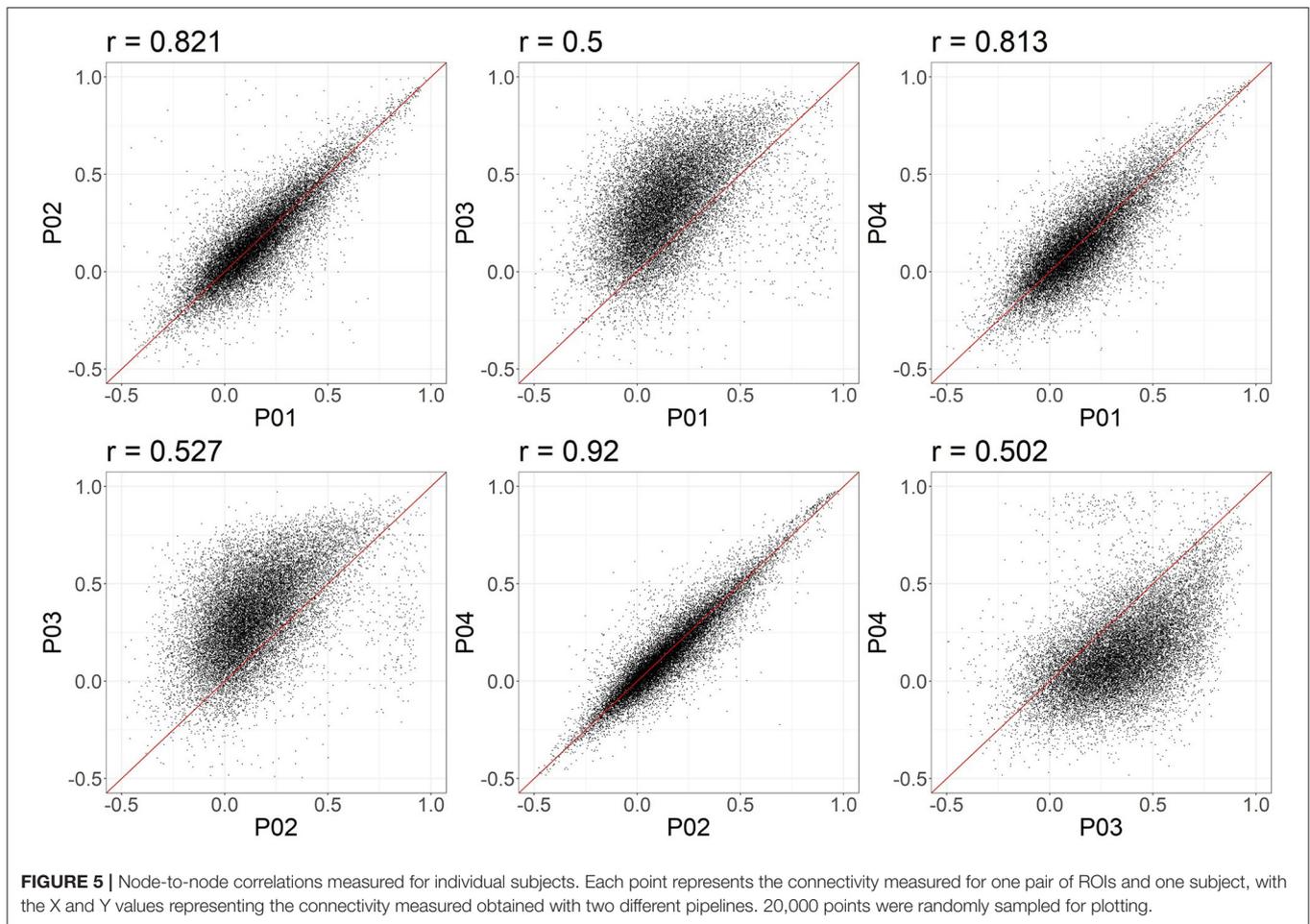
Along with the conversion of raw data into BIDS format, the Research Core also provides a set of analysis pipelines, training, and support.



Computational Environment

All data and analyses are hosted and completed on-site, providing full control of the systems' configuration and operation. Our specific implementation of the primary data storage is accomplished using a network attached storage cluster running the open-source Ceph file system (CephFS). We would like to note that any modern storage hardware/solution and/or mixed local storage with cloud storage should provide alternative option for another site implementation. We selected CephFS as a scalable solution installed on commodity hardware, which allows administrators to add storage incrementally without

rebuilding the entire cluster like some other solutions require. Performance scales with the size of the cluster, as data are not accessed through a fixed set of head nodes. The Laureate Institute for Brain Research (LIBR) currently has two petabytes of raw storage, which is 1PB of usable space after data duplication. Additionally, there is a full off-site backup copy stored roughly 100 miles away on an identical Ceph cluster. As a final precaution, LIBR also sends periodic tape backups to Iron Mountain using a Spectra BlackPearl appliance. Additional technical details relating to IT operations can be found in the supplement, Section 1.6.



LIBR has eight high-performance servers configured with the slurm workload manager (<https://slurm.schedmd.com/>). Each server has 24 physical cores, allowing up to 192 jobs to run in parallel and a total of over 24,000 GFlops/second. Jobs optimized to run on GPUs can take advantage of 4 Nvidia Tesla P100 cards, providing an additional 75,200 GFlops/second of computing power. Nodes are configured with 187 or 376 GB of RAM and overall networking throughput is 320 Gbps. This centralized processing infrastructure helps mitigate the bottleneck associated with network attached storage by providing 40 Gbps connections, which far outperform standard 1 Gbps connections used in modern ethernet.

The storage and computing infrastructure just described was designed and developed incrementally to balance cost with performance, security, and overhead for training and maintenance. As we noted above, our data organization and processing workflows, however, do not depend on the physical details of our environment and could be implemented on a variety of systems or in the cloud.

Pipeline Architecture Overview

Subject and group level analyses are conducted separately. Processing pipelines are implemented to service an individual

subject and analysis, where an analysis typically deals with one task and set of processing parameters (**Figure 2**). This allows for parallelization at the subject plus pipeline level, with separate jobs submitted for each subject.

All single-subject analyses are submitted to the batch scheduler using a script named preprocess-all-BIDS.py. This wrapper reads in a configuration file with pointers to the root of the source data directory (i.e., the root of one BIDS dataset), the desired root of the output directory tree and which pipeline to run. The output directory structure mirrors the BIDS formatted input, so that individual subject/session/pipeline results are stored in [Results Root]/sub-[subject]/ses-[session]/[leaf]. preprocess-all-BIDS.py traverses the input folder structure, and for every subject/session checks to see if a job has already been submitted, based on the existence of specially named status-indicating files in the output directory. If this subject/session combination has not been run for this pipeline, the output directory is created and a job is submitted.

Results for an individual subject/session/task/pipeline include derived values to be tabulated, quality control images in png or jpg format, and larger format derived data, like voxelwise statistics. Derived values include metrics like subject head motion, subject performance including mean reaction times and

accuracy, physiological measures including heart rate, and in the case of imaging tasks, extracted activations, contrasts, and volumes from atlas-based regions of interest. All derived values are stored in files ending in the .longformat suffix, where these are simple text files in attribute-value format. After processing data for all subjects, all values found in .longformat files are combined, producing a consolidated table with a single row per subject and session and one column per attribute. This consolidated format, ready for use in various statistics applications, is saved as in.csv and, RData formats, the later binary being preferable for large imaging datasets with tens of thousands of variables, which can lead to performance issues when reading in text data.

Any manual quality control processes are simplified by storing appropriate images in jpg or png format. For examples, this may include EKG traces with identified R wave peaks, or montages showing alignment and normalization of neuroimaging data. This allows the user to flip through QC images for a dataset relatively quickly without, for example, needing to open neuroimaging data in specialized software.

Neuroimaging Pipeline Options

fMRI Pipelines

Neuroimaging processing pipelines necessarily include numerous decisions, such as which software to use, whether to include linear or non-linear normalization to standard space, what smoothing kernel to apply, what nuisance regressors to use at the regression step and so on. These analysis decisions can impact the final results and interpretation of a study, which was recently illustrated through divergent results obtained by 70 independent groups of researchers who all analyzed the same data (15). Therefore, frameworks like ours that allow the sharing of analysis workflows are essential for reproducibility and replicability. An individual researcher may customize a particular pipeline or use one of our 3 standard options for each fMRI task, labeled P01 through P03. P01 is a traditional approach using AFNI (16) and includes removal of the first 3 volumes, despiking, slice-timing correction, co-registration between functional and structural volumes, motion correction, 4 mm of gaussian blur, and an affine transformation to standard space. P02 is similar to P01, except that it includes a non-linear warp to standard space and RETROICOR correction (17), which helps remove physiological noise but requires the collection of pulse oximeter and respiratory belt data.

P03 takes a completely different approach, instead using fMRIPrep (18) to do all preprocessing up until the regression step, which still uses AFNI's 3dDeconvolve. Preprocessing with fMRIPrep uses mainly default parameters, so that a combination of tools are used to (1) select a reference fMRI volume (mean of high contrast available in initial pre T1-saturation or pre Steady State Free Precession fMRI volume); (2) perform boundary based registration with the T1-weighted images (3, 19) estimate head motion prior to any spatiotemporal filtering using mcflirt in FSL 5.0.9 (4, 20) perform slice timing correction using AFNI (5, 16) perform nuisance regression including regressors for Framewise Displacement and DVARS (21); average CSE, white matter, and whole brain signals, as well as physiological regressors using CompCor (22). Regardless of the pipeline, standard derived

data from task-based fMRI include regression coefficients and contrasts extracted for each ROI in several atlases and summaries of head motion for quality control.

Resting state preprocessing P04 pipeline includes the same options as task data, with the addition of a fourth option, which is similar to P02 pipeline but also includes additional motion correction prior to slice timing correction via an automatic EEG assisted slice-specific motion correction for fMRI (aEREMCOR) (23). While it would be possible to include this additional motion correction step for task-based data, it is particularly important in resting state, where the residual effects of head motion are well known, and they might differ for each acquired slice (24). Standard derived data from resting-state fMRI include a correlation matrix between pairs of ROIs from multiple atlases [e.g., the Brainnetome (25)] and summaries of head motion.

EEG Pipelines

Simultaneous EEG-fMRI offers several benefits to measure and study the human brain's spatial and temporal dynamics in health and disease. However, EEG data collected during fMRI acquisition are contaminated with MRI gradients and ballistocardiogram artifacts, in addition to artifacts of physiological origin (eye blinks, muscle, motion), these artifacts need to be detected and suppressed before further data analysis (26). We have developed in house a comprehensive automated pipeline for EEG artifact reduction (APPEAR) recorded during fMRI, which we have incorporated into the BIDS preprocessing pipeline architecture (Figure 2). APPEAR is capable of reducing all main EEG artifacts, including MRI gradients, BCG, eye blinks, muscle, and motion artifacts, and can be applied to large (i.e., hundreds of subjects) EEG-fMRI datasets. APPEAR was evaluated, tested and compared to manual pre-processing EEG data for both resting EEG-fMRI recording as well as for event-related potential or task-based EEG-fMRI experiments in an exemplar eight subject EEG-fMRI dataset.

RESULTS

We provide examples illustrating pipelines P01 through P04 to help demonstrate the utility of the processing infrastructure. The full sample of participants who completed the NeuroMAP Core is summarized in **Supplementary Table 1**, and different subsets of these participants were used to provide example results from each pipeline.

Task fMRI Results

The Monetary Incentive Delay and Stop Signal tasks were included in NeuroMAP projects to probe the neural processing of reward and executive function (see **Supplementary Material** for task details). We have processed exemplar CDE data for each of these tasks and present the results in **Figure 3**. **Figure 3A** shows voxel-wise maps for the P5–P0 contrast in the MID as produced by pipelines P01 through P03. All three pipelines produce the expected activity in reward circuitry, however with some qualitative differences: For example, P01 produces the most widespread activation results, while P03 seems to identify more circumscribed areas. Data from 93 participants are included here,

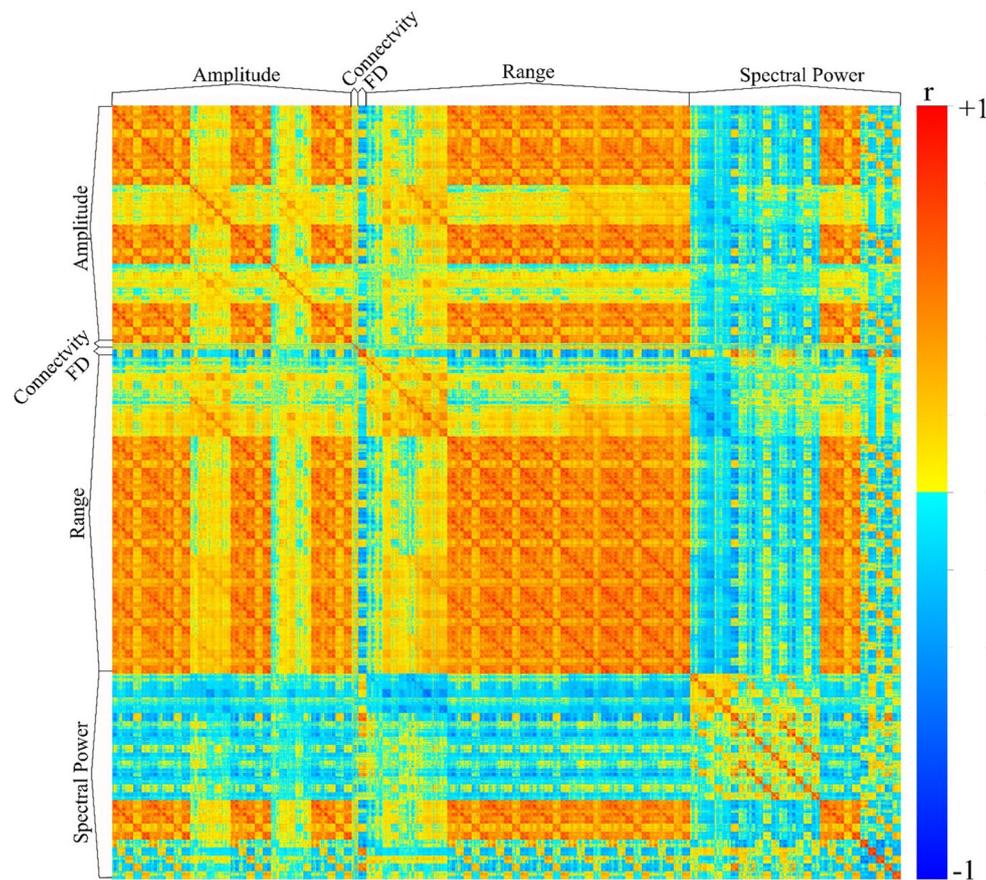


FIGURE 6 | The correlation matrix of 3,032 EEG features extracted using comprehensive EEG features extraction for resting-state condition. Five different subsets of features were extracted including Amplitude (31 Channels \times 5 bands \times 6 types = 930 features), connectivity (24 features), FD (31 Channels \times 1 Feature=31), range (31 Channels \times 5 bands \times 8 types = 1240 features) and spectral power features (31 Channels \times 5 bands \times 5 types + 31 Channels \times 1 Feature = 806 features). For more details about each subset of features, please see (28).

with pipelines P01 through P03 taking \sim 1.3, 4, and 5 CPU hours per subject to complete. With the architecture detailed in Section Pipeline Architecture Overview, processing for all three sets of data could be completed in under 1 day when all resources are available. The alignment QC images produced by each pipeline make it possible to complete all manual QC for roughly 100 participants and one pipeline in $<$ 1 h. **Figure 3B** shows voxelwise maps of the Stop–NoStop contrast from the stop signal task, again produced by pipelines P01 through P03. These maps include data from 49 participants. Again, maps show the expected activation (here primarily in parts of the executive control network), but with some qualitative differences. For example, P01 and P02 both identify negative activation in the motor cortex, which is not apparent in P03.

Resting State fMRI Results

Exemplar CDE data have also been processed for resting-state fMRI using all four pipelines, P01–P04. **Figure 4** shows the average connectivity matrix extracted from the Brainnetome atlas and organized by approximate networks identified using the Yeo 7-network atlas (27). All pipelines produce qualitatively similar

results at the group level, with functional networks apparent as colored squares on the diagonal.

In these data, and with the parameters we chose, P03 seems to identify consistently larger connectivity strengths compared to the other three pipelines. **Figure 5** shows the relationship between individual features (correlation strengths) measured with different pipelines. Points on the 45 degree line indicate complete agreement between methods, while divergence from that line illustrates differences between pipelines. This again illustrates the tendency for P03 to produce stronger correlations between pairs of ROIs. Of note, correlations between pairs of methods ranged from 0.5 to 0.9, meaning they only shared between approximately 25–80 percent variance. We cannot say which pipeline is best, but would highlight the large effect that selecting a particular pipeline may have on a the final results of a study.

EEG Preprocessing

We have utilized the APPEAR pipeline to preprocess EEG data acquired concurrently with fMRI, and then applied comprehensive EEG feature extraction from five subsets of EEG

features including amplitude, connectivity, fractal dimension (FD), range and spectral power features. Furthermore, each subset of features was applied to Alpha [8–13] Hz, Beta [15–30] Hz, Theta [4–7] Hz, Delta [0.5–4] Hz, Gamma [30–40] Hz and whole range of EEG frequency [0.5–40] Hz. An exemplar EEG feature correlation matrix is shown in **Figure 6**. The exemplar use of the extracted EEG features and automated EEG preprocessing can be found elsewhere [https://github.com/obadalzoubi/Comprehensive_EEG_Features_Extraction].

DISCUSSION

The proliferation of high-throughput data-generating technologies in biomedical research has led to data analytics challenges for creating easily reusable and reproducible pipelines. These challenges are especially salient for neuroscience studies, which not only involve the usual high-dimensional data but also include multiple neuroimage-specific data types and complex psychological trait data. The current study describes a scalable environment and set of software pipelines to preprocess neuroimaging (MRI, fMRI, and EEG) and behavioral data while integrating them with other subject-level high-dimensional data to perform sharable, reproducible analyses.

The services and computational environment developed by the Research Core provide a set of tangible benefits to ongoing research. Massive amounts of complex neuroimaging data are put into a standard (BIDS) format with minimal human interaction in an ongoing basis. The architecture for converting data to BIDS format is flexible and scalable, so that new studies often have compliant data from day 1.

Once the data for a study are in BIDS format, running any of our standard preprocessing pipelines becomes a quick process. With relatively little human intervention, preprocessing jobs can be created for hundreds or thousands of participants, and the processing and network storage infrastructure can produce results in days rather than weeks. Having multiple pipelines available for the same tasks gives researchers the ability to verify that their results are robust to the details of the preprocessing pipeline, as others have shown the wide variation in analysis results to be a serious concern (15).

In this work, we provide exemplar results 11 different pipelines (three pipelines on each of two fMRI tasks, four pipelines on resting-state fMRI, and one pipeline on resting EEG) to demonstrate the utility of our infrastructure. Additionally, ROI-level results from our standard pipelines have been used in studies of cannabis (29) and stimulant/opioid use (30), while voxelwise results have appeared in studies of neighborhood effects (31) and inflammation (32), and clinical data have been used to predict head motion during scanning (33). We have also used EEG derived features have to differentiate participants with mood and anxiety disorders from healthy controls (34) and to predict participant age (28).

Our workflow incorporates many diverse processing and analysis tools such as *afni*, *freesurfer*, *fmrip* and uses the BIDS format. However, it has been noted that the large number of analysis degrees of freedom in neuroscience increases the risk of false discoveries due *p*-value hacking or “researcher degrees of freedom” (35). Each analysis step can result in an expanding decision tree of potential analyses. Determining the best workflow software or pipeline option for a given experiment is an ongoing question, but the current software provides standard selections for the many analysis options. As the field evolves and standards consolidate, the default processing and analysis parameters will converge to standards with lower variation and increased replicability.

In addition to neuroimaging data, our current pipelines include other common data types and can be easily extended to other high-throughput data, such as genetic and gene expression. Many neuroscience studies also include large non-neuroimage datasets, such as GWAS, which has its own relatively complex file format known as Plink. BIDS is open source and under active development, and integration with these other datasets will be straightforward extensions of BIDS.

We hope the data management and processing infrastructure presented here may act as a blueprint for other organizations seeking to standardize data collection and processing. We also hope this serves as a testament to the growing widespread adoption of BIDS as a common standard. Some of our tools are publicly available on GitHub (where noted above), and others may be shared on reasonable request.

ETHICS STATEMENT

The studies involving human participants were reviewed and approved by Western IRB. The patients/participants provided their written informed consent to participate in this study and received financial compensation for participation.

AUTHOR CONTRIBUTIONS

MP, JB, BM, RA, and TT contributed to conception and design of the study. RK, MM, OA, AM, JB, and JT contributed to analytic workflows development. RK and OA performed the analyses. RK and JB drafted the manuscript. JT, OA, RK, BM, AM, MM, and JB wrote sections of the manuscript. All authors contributed to manuscript revisions, read, and approved the submitted version.

FUNDING

This work was supported by National Institute of General Medical Sciences, National Institutes of Health P20GM121312 award, and in part by in part by W81XWH-12-1-0697 award from the U.S. Department of Defense, the Laureate Institute for Brain Research (LIBR), and the William K. Warren Foundation. The funding agencies

were not involved in the design and development, data collection and analyses, and preparation and submission of the manuscript.

ACKNOWLEDGMENTS

We thank Dr. Jennifer Stewart for valuable training contributions for NeuroMAP-Investigators. The NeuroMAP-Investigators include the following contributors: Yoon-Hee Cha MD., Justin S.

Feinstein Ph.D., Sahib S. Khalsa MD., Jonathan Savitz Ph.D., W. Kyle Simmons Ph.D., Namik Kiric Ph.D., Maria Ironside Ph.D., and Evan J. White Ph.D.

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fpsy.2021.682495/full#supplementary-material>

REFERENCES

- Leow AD, Yanovsky I, Parikshak N, Hua X, Lee S, Toga AW, et al. Alzheimer's disease neuroimaging initiative: a one-year follow up study using tensor-based morphometry correlating degenerative rates, biomarkers and cognition. *Neuroimage*. (2009) 45:645–55. doi: 10.1016/j.neuroimage.2009.01.004
- 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. doi: 10.1016/j.neuroimage.2013.05.041
- Jernigan TL, Brown SA, Dowling GJ. The Adolescent Brain Cognitive Development Study. *J Res Adolesc*. (2018) 28:154–6. doi: 10.1111/jora.12374
- Victor TA, Khalsa SS, Simmons WK, Feinstein JS, Savitz J, Aupperle RL, et al. Tulsa 1000: a naturalistic study protocol for multilevel assessment and outcome prediction in a large psychiatric sample. *BMJ Open*. (2018) 8:e016620. doi: 10.1136/bmjopen-2017-016620
- Gorgolewski KJ, Auer T, Calhoun VD, Craddock RC, Das S, Duff EP, et al. The brain imaging data structure, a format for organizing and describing outputs of neuroimaging experiments. *Sci Data*. (2016) 3:160044. doi: 10.1038/sdata.2016.44
- Marcus DS, Olsen TR, Ramaratnam M, Buckner RL. The Extensible Neuroimaging Archive Toolkit: an informatics platform for managing, exploring, and sharing neuroimaging data. *Neuroinformatics*. (2007) 5:11–34. doi: 10.1385/ni:5:1:11
- Keator DB, Grethe JS, Marcus D, Ozyurt B, Gadde S, Murphy S, et al. A national human neuroimaging collaborative enabled by the Biomedical Informatics Research Network (BIRN). *IEEE Trans Inf Technol Biomed*. (2008) 12:162–72. doi: 10.1109/TITB.2008.917893
- Van Horn JD, Toga AW. Is it time to re-prioritize neuroimaging databases and digital repositories? *Neuroimage*. (2009) 47:1720–34s. doi: 10.1016/j.neuroimage.2009.03.086
- Ozyurt IB, Keator DB, Wei D, Fennema-Notestine C, Pease KR, Bockholt J, et al. Federated web-accessible clinical data management within an extensible neuroimaging database. *Neuroinformatics*. (2010) 8:231–49. doi: 10.1007/s12021-010-9078-6
- Das S, Zijdenbos AP, Harlap J, Vins D, Evans AC. LORIS: a web-based data management system for multi-center studies. *Front Neuroinform*. (2011) 5:37. doi: 10.3389/fninf.2011.00037
- Scott A, Courtney W, Wood D, de la Garza R, Lane S, King M, et al. COINS: an innovative informatics and neuroimaging tool suite built for large heterogeneous datasets. *Front Neuroinform*. (2011) 5:33. doi: 10.3389/fninf.2011.00033
- Book GA, Anderson BM, Stevens MC, Glahn DC, Assaf M, Pearlson GD. Neuroinformatics Database (NiDB)—a modular, portable database for the storage, analysis, and sharing of neuroimaging data. *Neuroinformatics*. (2013) 11:495–505. doi: 10.1007/s12021-013-9194-1
- Harris PA, Taylor R, Thielke R, Payne J, Gonzalez N, Conde JG. Research electronic data capture (REDCap)—a metadata-driven methodology and workflow process for providing translational research informatics support. *J Biomed Inform*. (2009) 42:377–81. doi: 10.1016/j.jbi.2008.08.010
- Harris PA, Taylor R, Minor BL, Elliott V, Fernandez M, O'Neal L, et al. The REDCap consortium: Building an international community of software platform partners. *J Biomed Inform*. (2019) 95:103208. doi: 10.1016/j.jbi.2019.103208
- Botvinik-Nezer R, Holzmeister F, Camerer CF, Dreber A, Huber J, Johannesson M, et al. Variability in the analysis of a single neuroimaging dataset by many teams. *Nature*. (2020) 582:84–88. doi: 10.1038/s41586-020-2314-9
- Cox RW, Hyde JS. Software tools for analysis and visualization of fMRI data. *NMR Biomed*. (1997) 10:171–8. doi: 10.1002/(sici)1099-1492(199706/08)10:4/5<171::aid-nbm453>3.0.co;2-1
- Glover GH, Li TQ, Ress D. Image-based method for retrospective correction of physiological motion effects in fMRI: RETROICOR. *Magn Reson Med*. (2000) 44:162–7. doi: 10.1002/1522-2594(200007)44:1<162::aid-mrm23>3.0.co;2-e
- Esteban O, Markiewicz CJ, Blair RW, Moodie CA, Isik AI, Erramuzpe A, et al. fMRIPrep: a robust preprocessing pipeline for functional MRI. *Nature methods*. (2019) 16:111–6. doi: 10.1038/s41592-018-0235-4
- Greve DN, Fischl B. Accurate and robust brain image alignment using boundary-based registration. *Neuroimage*. (2009) 48:63–72. doi: 10.1016/j.neuroimage.2009.06.060
- Jenkinson M, Bannister P, Brady M, Smith S. Improved optimization for the robust and accurate linear registration and motion correction of brain images. *Neuroimage*. (2002) 17:825–41. doi: 10.1016/s1053-8119(02)91132-8
- Power JD, Mitra A, Laumann TO, Snyder AZ, Schlaggar BL, Petersen SE. Methods to detect, characterize, and remove motion artifact in resting state fMRI. *Neuroimage*. (2014) 84:320–41. doi: 10.1016/j.neuroimage.2013.08.048
- Behzadi Y, Restom K, Liao J, Liu TT. A component based noise correction method (CompCor) for BOLD and perfusion based fMRI. *Neuroimage*. (2007) 37:90–101. doi: 10.1016/j.neuroimage.2007.04.042
- Wong, C.-K., Zotev V, Misaki M, Phillips R, Luo Q, et al. Automatic EEG-assisted retrospective motion correction for fMRI (aE-REMCOR). *Neuroimage*. (2016) 129:133–47. doi: 10.1016/j.neuroimage.2016.01.042
- Power JD, Schlaggar BL, Petersen SE. Recent progress and outstanding issues in motion correction in resting state fMRI. *Neuroimage*. (2015) 105:536–51. doi: 10.1016/j.neuroimage.2014.10.044
- Fan L, Li H, Zhuo J, Zhang Y, Wang J, Chen L, et al. The Human brainnetome atlas: a new brain atlas based on connectational architecture. *Cereb Cortex*. (2016) 26:3508–26. doi: 10.1093/cercor/bhw157
- Mayeli A, Henry K, Wong CK, Zoubi OA, White EJ, Luo Q, et al. (2019) Automated Pipeline for EEG Artifact Reduction (APPEAR) recorded during fMRI. *arXiv [Preprint]* arXiv:1912.05507. Available online at: <https://arxiv.org/abs/1912.05507>
- Yeo BT, Krienen FM, Sepulcre J, Sabuncu MR, Lashkari D, Hollinshead M, et al. The organization of the human cerebral cortex estimated by intrinsic functional connectivity. *J Neurophysiol*. (2011) 106:1125–65. doi: 10.1152/jn.00338.2011
- Al Zoubi O, Ki Wong C, Kuplicki RT, Yeh HW, Mayeli A, Refai H, et al. Predicting Age From Brain EEG signals—a machine learning approach. *Front Aging Neurosci*. (2018) 10:184. doi: 10.3389/fnagi.2018.00184
- Spechler PA, Stewart JL, Kuplicki R, Tulska I, Paulus MP. Attenuated reward activations associated with cannabis use in anxious/depressed individuals. *Transl Psychiatry*. (2020) 10:189. doi: 10.1038/s41398-020-0807-9
- Stewart JL, Khalsa SS, Kuplicki R, Puhl M, Investigators T, Paulus MP. Interoceptive attention in opioid and stimulant use disorder. *Addict Biol*. (2020) 25:e12831. doi: 10.1111/adb.12831

31. Feng C, Forthman KL, Kuplicki R, Yeh HW, Stewart JL, Paulus MP. Neighborhood affluence is not associated with positive and negative valence processing in adults with mood and anxiety disorders: A Bayesian inference approach. *Neuroimage Clin.* (2019) 22:101738. doi: 10.1016/j.nicl.2019.101738
32. Burrows K, Stewart JL, Kuplicki R, Figueroa-Hall L, Spechler PA, Zheng H, et al. Elevated peripheral inflammation is associated with attenuated striatal reward anticipation in major depressive disorder. *Brain Behav Immunity.* (2021) 93:214–25. doi: 10.1016/j.bbi.2021.01.016
33. Ekhtiari H, Kuplicki R, Yeh HW, Paulus MP. Physical characteristics not psychological state or trait characteristics predict motion during resting state fMRI. *Sci Rep.* (2019) 9:419. doi: 10.1038/s41598-018-36699-0
34. Al Zoubi O, Mayeli A, Tsuchiyagaito A, Misaki M, Zotev V, Refai H, et al. EEG microstates temporal dynamics differentiate individuals with mood and anxiety disorders from healthy subjects. *Front Hum Neurosci.* (2019) 13:56. doi: 10.3389/fnhum.2019.00056
35. Wicherts JM, Veldkamp CL, Augusteijn HE, Bakker M, van Aert RC, van Assen MA. Degrees of freedom in planning, running, analyzing, and reporting psychological studies: a checklist to avoid p-hacking. *Front Psychol.* (2016) 7:1832. doi: 10.3389/fpsyg.2016.01832

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2021 Kuplicki, Touthang, Al Zoubi, Mayeli, Misaki, NeuroMAP-Investigators, Aupperle, Teague, McKinney, Paulus and Bodurka. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.