

Interpreting models interpreting brain dynamics

Md Mahfuzur Rahman (✉ mahfuz.gsu@gmail.com)

Tri-institutional Center for Translational Research in Neuroimaging and Data Science

<https://orcid.org/0000-0002-4162-6152>

Usman Mahmood

Tri-institutional Center for Translational Research in Neuroimaging and Data Science

Noah Lewis

Tri-institutional Center for Translational Research in Neuroimaging and Data Science

Harshvardhan Gazula

Massachusetts General Hospital and Harvard Medical School

Alex Fedorov

Tri-institutional Center for Translational Research in Neuroimaging and Data Science

Zening Fu

Tri-Institutional Center for Translational Research in Neuroimaging and Data Science (TReNDS)

<https://orcid.org/0000-0002-1591-4900>

Vince Calhoun

Tri-Institutional Center for Translational Research in Neuroimaging and Data Science (TReNDS): Georgia State University Georgia Institute of Technology and Emory University <https://orcid.org/0000-0001-9058-0747>

Sergey Plis

Tri-institutional Center for Translational Research in Neuroimaging and Data Science

Article

Keywords: brain dynamics, deep learning, resting-state fMRI

Posted Date: December 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-798060/v2>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Interpreting models interpreting brain dynamics

Md. Mahfuzur Rahman^{1, 2, *}, Usman Mahmood^{1, 2}, Noah Lewis^{1, 3}, Harshvardhan Gazula⁴, Alex Fedorov^{1, 5}, Zening Fu¹, Vince D. Calhoun^{1, 2, 3, 5}, and Sergey M. Plis^{1, 2}

¹Tri-institutional Center for Translational Research in Neuroimaging and Data Science (TReNDS), Georgia State University, Georgia Institute of Technology, Emory University, Atlanta, GA, USA

²Georgia State University, Department of Computer Science, Atlanta, GA, USA

³Georgia Institute of Technology, School of Computational Science & Engineering, Atlanta, GA, USA

⁴Athinoula A. Martinos Center for Biomedical Imaging, Massachusetts General Hospital and Harvard Medical School, MA, USA

⁵Georgia Institute of Technology, School of Electrical & Computer Engineering, Atlanta, GA, USA

*mahfuz.gsu@gmail.com

ABSTRACT

Brain dynamics are highly complex and yet hold the key to understanding brain function and dysfunction. The dynamics captured by resting-state functional magnetic resonance imaging data are noisy, high-dimensional, and not readily interpretable. The typical approach of reducing this data to low-dimensional features and focusing on the most predictive features comes with strong assumptions and can miss essential aspects of the underlying dynamics. In contrast, introspection of discriminatively trained deep learning models may uncover disorder-relevant elements of the signal at the level of individual time points and spatial locations. Yet, the difficulty of reliable training on high-dimensional low sample size datasets and the unclear relevance of the resulting predictive markers prevent the widespread use of deep learning in functional neuroimaging. In this work, we introduce a deep learning framework to learn from high-dimensional dynamical data while maintaining stable, ecologically valid interpretations. Results successfully demonstrate that the proposed framework enables learning the dynamics of resting-state fMRI directly from small data and capturing compact, stable interpretations of features predictive of function and dysfunction.

Introduction

Brain dynamics likely holds the key to understanding function and disorder¹⁻³. The brain function manifests in a spatiotemporally localized activity within the dynamics⁴. Thus, identification and interpretation of subject-specific spatial and temporal activity may help guide our understanding of the disorder. Although, the spatiotemporal snapshots of brain dynamics can be captured noninvasively using functional magnetic resonance imaging (fMRI)^{5,6}, the excessive dimensionality and complexity of fMRI signals rule out manual identification and interpretation. Alternatively, machine learning models trained to classify a mental disorder from the available observations have learned which aspects of the data reliably lead to correct prediction. In other words, the model builds internal representations of the mapping between the data and the class. Interpreting these representations can lead to discovery of previously unknown spatiotemporal functional indicators (or biomarkers).

However, standard machine learning (SML) models, when dealing directly with high-dimensional multivariate signals, suffer a drastic drop in performance because of the curse of dimensionality⁷ (high dimensionality of fMRI relative to the typically available few samples). To deal with this issue, neuroimaging researchers resort to hand-engineered features⁸ such as correlation matrices, also called Functional Network Connectivity (FNC), that summarize spatiotemporal relationship between different brain regions^{9,10}. Arguably, such proxy representations rely on strict assumptions and miss the chance to discover highly predictive holistic representations of the underlying dynamics^{11,12}. This limitation calls for a shift from a feature-engineering paradigm to a feature-learning paradigm that can allow model introspection and automatic discovery of the spatiotemporal activity indicative of the disorder under consideration. When available, such a feature learning paradigm may greatly facilitate discovering actionable causal knowledge about the disorder.

In recent years, deep learning (DL) models have attracted significant attention for their ability to learn reliable and robust features directly from the high-dimensional data in diverse neuroimaging applications¹³⁻¹⁶ in addition to their highly discriminative capabilities. More recently, Abrol et al. (2021)¹⁷ demonstrated the advantages of DL models trained on raw data over SML models trained on pre-engineered features in structural magnetic resonance imaging (sMRI). They also demonstrated the potential to interpret and visualize discriminatory biomarkers within the data by leveraging robust introspection techniques. The study suggests that the deep representations of dynamics (fMRI) may be as discriminative and informative as their structural counterparts (sMRI). However, not every DL model is simultaneously predictive and interpretable for time series data capturing

39 dynamics¹⁸.

40 The predictive performance of a DL model is strongly proportional to the size of training data¹⁹, which in most neuroimaging
41 studies is scarce to come by due to the costly data collection process. In such a scenario, transfer learning can be a convenient
42 approach to dealing with this problem, as reported in numerous studies²⁰⁻²⁴. Although transfer learning usually involves
43 supervised pretraining of a model on a related task, it is difficult to find a way to formulate the pretraining task and also the data
44 to use so as to benefit the downstream fMRI tasks. Model interpretation may be challenging for overparameterized models, but
45 if the architecture supports robust and stable sensitivity analyses^{25,26}, the interpretations for individual predictions will also be
46 stable and robust.

47 The main idea of this paper is that DL can learn directly from high-dimensional signal dynamics even in small datasets
48 and, upon introspection, can help discover disease-specific salient data regions, which, if carefully utilized, can advance our
49 understanding of brain function. To achieve this, we introduce a model that learns from dynamical data and lends itself to
50 interpretations. To maximally benefit from small data, we propose a self-supervised pretraining scheme^{22,23}, which maximises
51 “mutual information local to (whole) context” *whole MILC*, to capture potentially valuable knowledge from the data not directly
52 related to the study. Our pretraining leverages publicly available healthy control subjects from the Human Connectome Project
53 (HCP)²⁷ to establish prior knowledge about the general signal dynamics and directly transfer the insights into the downstream
54 small data studies of schizophrenia, autism, and Alzheimer’s disease with subject age-range significantly broader than in
55 HCP. Subsequently, to validate the discovered biomarkers, we propose a “Retain And Retrain” (*RAR*) method to show that
56 the biomarkers identified as explanations are demonstrably informative. In particular, *RAR* equipped with an SML model
57 can verify and quantify the effectiveness of the feature attributions. More precisely, the identified salient features are highly
58 predictive compared to random baselines and, as we further show, capture the essence of the disorder-specific brain dynamics.
59 A visual depiction¹ of the proposed framework is shown in Figure 1.

61 Results

62 We first describe all the datasets and present the results under two broad sections—*whole MILC Performance* and *Post hoc*
63 *Explanation & RAR Evaluation on FNC*. The *whole MILC performance* indicates its predictive capacity in discriminating
64 patients from healthy controls for each disorder separately. *Post hoc explanations* are feature attributions as determined by the
65 *whole MILC* model for its predictions which we subsequently evaluated using the *RAR* scheme via an independent SVM
66 model.

67 Datasets

68 We used the Autism Brain Imaging Data Exchange (ABIDE)²⁸ (569 subjects- 255 healthy controls (HC) and 314 patients)
69 for autism spectrum disorder (ASD), the Function Biomedical Informatics Research Network (FBIRN)²⁹ (311 subjects- 151
70 healthy controls and 160 patients) for schizophrenia (SZ), and the Open Access Series of Imaging Studies (OASIS)³⁰ (372
71 subjects- 186 healthy controls and 186 patients) for Alzheimer’s disease (AZ).

72 whole MILC Performance

73 We evaluated the effectiveness of the proposed DL architecture with (w/) and without (w/o) the proposed self-supervised
74 pretraining scheme, aka *whole MILC*, by comparing its performance against standard machine learning models. We also
75 progressively increased the downstream sample size to investigate its impact on the model’s discriminative capacity. We used a
76 K-fold cross-validation strategy for all the experiments below. The model was trained on samples progressively selected from
77 the train folds, and we report the performance (AUC) on the test fold.

78 whole MILC Evaluation

79 **Autism (ABIDE) Results** (with K = 6) (see Figure 2 Autism spectrum panel) show that when we used a small number of
80 subjects for training (e.g., 15 subjects per class), the pretraining improved the model’s performance compared to when the
81 model learned only from the downstream training data (“w/o pretraining”). However, as we gradually increased the number
82 of training subjects, the model w/o pretraining outperformed the model w/ pretraining. The reduced effects of pretraining on
83 autism disorder classification are reasonable because the subjects from the HCP dataset are from different age groups than
84 those from the ABIDE dataset.

85 **Schizophrenia (FBIRN) Results** (with K = 5) (see Figure 2 Schizophrenia panel) show that the proposed architecture w/
86 pretraining outperformed w/o pretraining at almost all sample sizes, and the difference was more pronounced at smaller sample
87 sizes.

¹Human silhouettes in Figure 1 are by Natasha Sinegina for Creazilla.com without modifications, <https://creativecommons.org/licenses/by/4.0/>

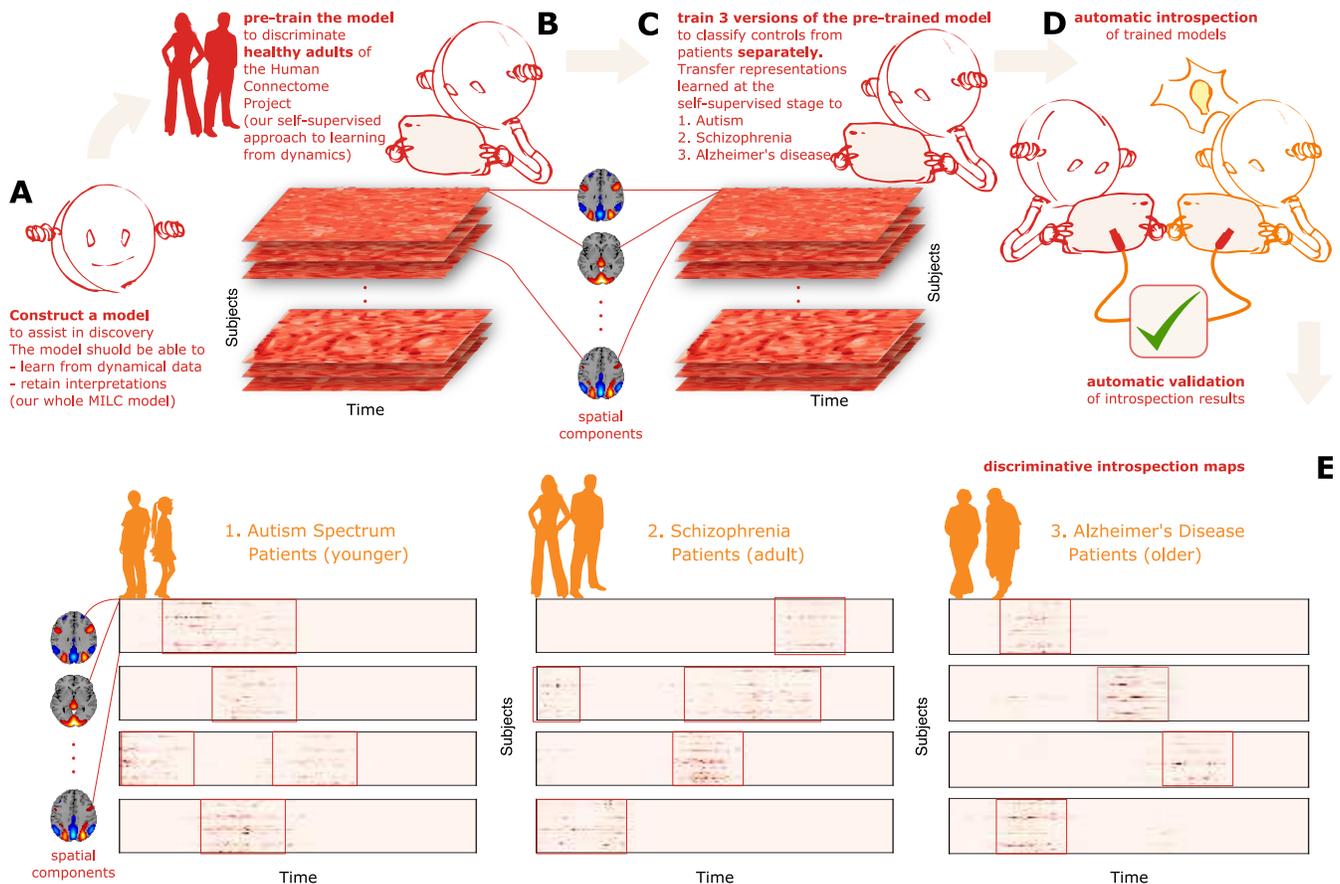


Figure 1. An overview of our approach to model interpretation. **A:** Construct a model for disorder-specific discovery: the *whole MILC* model learns directly from the disorder signal dynamics and retains interpretations for further introspection. **B:** Leverage self-supervised pretraining to distinguish healthy subjects: learned representations assist the model in maintaining its predictive power when downstream training data is limited. **C:** Construct a downstream model to discriminate patients from controls for each disorder starting with the pre-trained *whole MILC* weights: transfer of representations learned during pretraining simplifies convergence and balances overfitting. **D:** Introspection of the trained downstream models: interpretability methods extract meaningful, distinctive parts through feature attributions. Subsequently, the estimated salient aspects of the dynamics go through an automatic validation process. To this end, we use the most salient features to retrain an independent SML model that confirms the salience of the features. This information can then be relayed to a human expert in the relevant field to interpret further and advance knowledge about the disorders. **E:** Examples of saliency maps as deemed highly predictive by the models for their predictions in three different discriminative tasks.

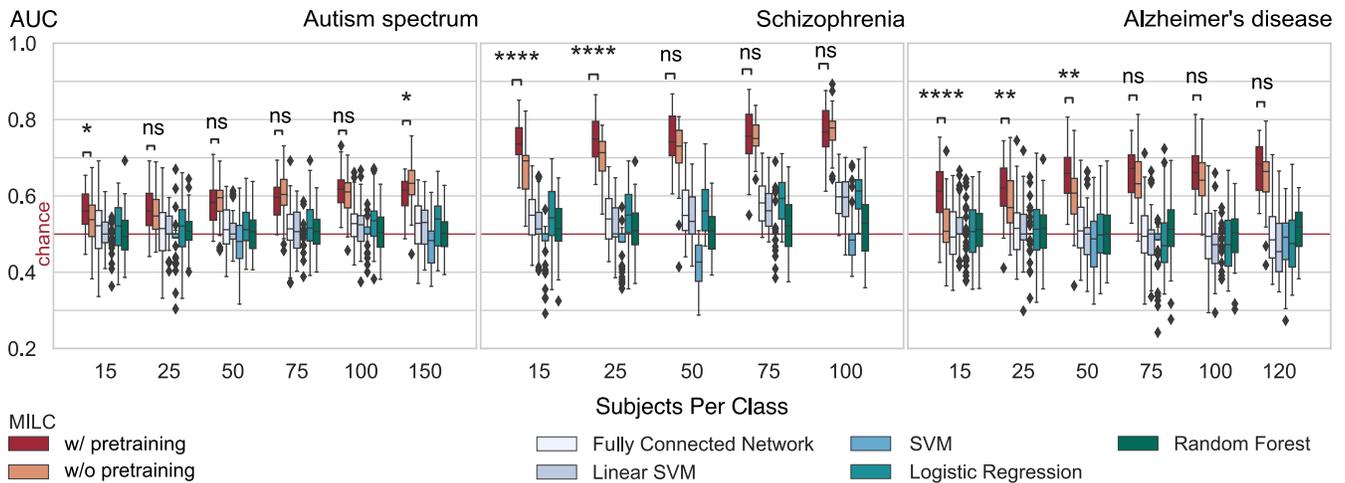


Figure 2. The main results from the whole MILC architecture and its comparison with standard machine learning models (SML). Apparently, the whole MILC model, in general, can learn from the raw data where traditional SML models fail to maintain their predictive capacity. Moreover, the whole MILC w/ pretraining substantially improves the latent representations as reflected in the improved accuracy compared to the whole MILC w/o pretraining. Specifically, in most small data cases, the whole MILC w/ pretraining outperformed the whole MILC w/o pretraining across the datasets. However, as expected, when we gradually increased the number of subjects during training, the effect of pretraining on the classification performance diminished, and both configurations of whole MILC did equally well. We verified this trend over three datasets that correspond to autism spectrum disorder, schizophrenia, and Alzheimer’s disease.

88 **Alzheimer’s disease (OASIS)** Similar to what has been observed in the case of SZ (FBIRN), the effect of pretraining on
 89 the downstream classification task ($K = 6$) (see Figure 2 Alzheimer’s disease panel) was more pronounced (comfortably
 90 outperforming) than w/o pretraining. This margin was substantial when the training data size was limited. However, as we
 91 increased the training data size, the gap between "w/ pretraining" and "w/o pretraining" was hardly conceivable.

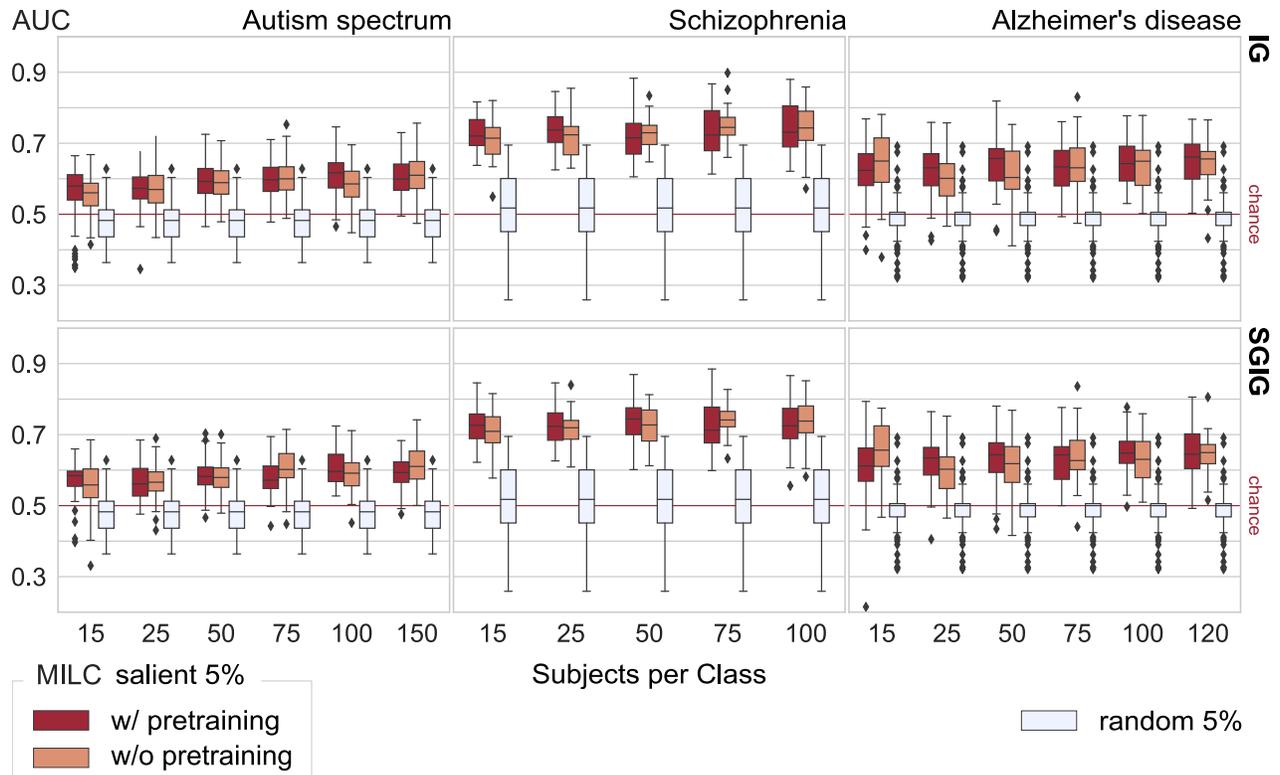
92 **Post hoc Explanation & RAR Evaluation using FNC**

93 Once the *whole MILC* model was trained, we computed the feature attributions (saliency maps) as determined by the model for
 94 each prediction. These feature attribution values were estimated for every subject from the dataset because the subsequent
 95 validation depends on training and test samples. Using RAR and an independent SVM model, we validated the model-identified
 96 salient parts of data to demonstrate that the highly regarded input parts were empirically discriminative and meaningful. Before
 97 RAR evaluation, we computed the average importance values of the overlapped time steps to obtain a single attribution value
 98 for every spatiotemporal dimension in the input sample. Refer to Figure 1 for example maps of patients from all the relevant
 99 disorder datasets.

100 **RAR Evaluation**

101 For RAR evaluation, we trained an SVM model on FNC matrices measured as Pearson’s correlations between time courses of
 102 the components obtained by spatial independent component analysis (ICA)³¹ (discussed in Methods section). We estimated
 103 this FNC based on only 5% salient or random (baseline) data. The RAR validation results of different models trained on three
 104 datasets with the most salient 5% (see Supplementary Fig. 1 for results from different percentages of salient data) training
 105 data are reported in Figure 3. As we can see, the dynamics learned by the *whole MILC* model were essential to maintain
 106 its predictive capacity. We observed that the model-specified salient data parts were more predictive than a similar amount
 107 of randomly chosen input data when we evaluated them for the same classification task using an independent SVM. This
 108 encouraging performance based on the salient data implies that the model can capture spatiotemporally meaningful markers
 109 suitable for patient-control distinction. Moreover, in many cases, the biomarkers identified with the "w/ pretraining"
 110 variant of the *whole MILC* model were more discriminative than the biomarkers specified with the "w/o pretraining"
 111 version, as reflected in the SVM’s classification performance. This encouraging result generalized across the datasets, even when we used very few
 112 subjects (15) for training.

113 As demonstrated in classification performance shown in Figure 2 and validation of feature attributions shown in Figure
 114 3, it is evident that the three predictive tasks were successful using our transfer learning model. In addition to quantitative
 115 validation of the automatic model introspection, we further analyzed the group-level functional network connectivity based



Reference single-subject saliency mask at varying degrees of data coverage



Figure 3. RAR employs SVM to classify the FNCs of the top 5% of the salient input data as estimated by the *whole MILC* model's predictions. We used integrated gradients (IG) and smoothgrad integrated gradients (SGIG) to compute feature attributions. It is evident that when an independent classifier (SVM) learned on every subject's most salient 5% data, the predictive power was significantly higher compared to the same SVM model trained on the randomly chosen same amount of data. In other words, the poor performance with randomly selected data parts indicates that other parts of the data were not exclusively discriminative as the *whole MILC* estimated salient 5% data parts. We also notice that sample masks over a different percentage of data coverage gradually obscured the localization of the discriminative activity within the data. Though the SVM model gradually became predictive with increased randomly selected data coverage, which we show in Supplementary Information, this performance upgrade was due to the gradual improvement in functional connectivity estimation and not attributable to the disease-specific localized parts within the data. For every disorder (Autism spectrum disorder, Schizophrenia, and Alzheimer's disease), the higher AUC at this 5% indicates stronger relevance of the salient data parts to the underlying disorders. Furthermore, the RAR results reflect that in most cases, when *whole MILC* was trained with limited data, the w/ pretraining models estimated feature attributions more accurately than the models w/o pretraining.

116 on the model-identified salient parts of data. Refer to the connectograms (see Figure 4) showing the top 10% FNC computed
117 using the most 5% discriminative data as localized by the trained model for the patients in three different disorders. We can see
118 some interesting differences in the connectograms. Autism spectrum disorder (ABIDE) shows the least between-domain FNC
119 highlighting within domain changes in specific cerebellum, sensorimotor, and subcortical domains³². Schizophrenia (FBIRN)
120 has the most widespread predictive pattern, consistent with prior work³³ showing cerebellum interaction across multiple
121 domains and sensorimotor changes. Finally, the predictive features for Alzheimer’s disease (OASIS) are mainly engaging visual
122 and cognitive interactions³⁴. Figure 5 shows full FNC matrices (based on 5% data), their disorder pairwise difference, and
123 static FNC matrices (based on 100% data) for all disorders. As we can observe, the proposed model could capture the essential
124 dynamics as generally captured in traditional full data FNC matrices and thus fully consistent with the knowledge from existing
125 literature. The pairwise difference matrices imply that the different brain dynamics are indeed different for different disorders.

126 Furthermore, we also investigated the temporal characteristics of the saliency maps for patients and controls of each disorder.
127 For this, we first determined the most important time points for each saliency map, expressed as temporal density and computed
128 as the number of components for each time point that appeared in the top 5% values of the map. We observed interesting
129 differences between groups in temporal behavior. In particular, we noticed that the temporal behavior of the most discriminative
130 time steps is much more focused for schizophrenia and Alzheimer’s patients than their healthy controls counterparts. Put
131 another way, the temporal density of schizophrenia and Alzheimer’s patients is generally spiky, whereas, for the healthy controls
132 it is largely flatter. However, for autism spectrum disorder, the temporal density behavior between patients and controls is
133 largely uniform, and the distinction, if any, is hardly noticeable. Refer to Figure 6 panel A for some samples showing temporal
134 behavior of patients and controls for all disorders. To quantify these temporal characteristics (spikiness and uniformity in
135 temporal densities), we calculated the earth mover’s distance (EMD)³⁵—a distance measure between two densities—between
136 the temporal density computed from each saliency map and a uniform density function. The intuition behind this spread
137 measure is that a small EMD indicates that the distribution is predominantly uniform and not localized in time, implying that
138 the discriminatory activity is usually not confined to any specific time interval. On the other side, a large EMD indicates
139 spikiness of the temporal behavior signaling that the discriminative activity is more focused in a shorter time interval. Refer
140 to Figure 6 panel B for the distributions of EMD and corresponding statistical test results for all the disorders. We observe
141 that the discriminative activity for schizophrenia patients is predominantly local and hence more focused in time, whereas the
142 distinguishing characteristics of healthy controls are spread across time. We observed similar characteristics for Alzheimer’s
143 patients. However, for autism spectrum disorder, we noticed that the temporal characteristics for both patients and controls are
144 generally spread across time and not distinguishable. We verified our observations through a non-parametric statistical test
145 conducted on EMD distributions for each disorder.

146 Discussion

147 Standard machine learning models are widely used in neuroimaging research partly due to their familiarity and ease of use
148 and the perceived simplicity of interpretability of the outcomes. However, this ease/simplicity takes a hit when the complexity
149 and dimensionality of the input data are high, as is often the case with fMRI data. Our experiments (Figure 2) show that SML
150 models fail to achieve good predictive performance, let alone provide meaningful interpretations of the underlying dynamics.
151 This failure is not surprising since these proxy features are sensitive to strict assumptions about the signal dynamics^{11,12}, which
152 may only be partially accurate or accurate just under certain conditions. However, deep learning models can overcome this
153 curse of dimensionality and learn meaningful interpretations in addition to showing high predictive performance^{15–17}. This
154 work demonstrates that DL models can achieve a deeper understanding of the underlying subject-specific signal dynamics in an
155 fMRI setting despite the commonly expected difficulty of interpretability.

156 While recent advances in deep learning have proved its impressive ability to learn from a signal close to the raw data,
157 different network architectures have benefits and limitations. The default choice of deep learning architecture for time-series
158 data is the well-known RNN class of models, specifically LSTM. Although LSTM models return good performance, they still
159 have issues with interpretability due to vanishing saliency, making them unsuitable for studying multivariate signal dynamics.
160 This necessitates building a suitable architecture that can resolve the vanishing saliency problem in the recurrent model while
161 preserving the stability and making attributions meaningful. To that end, Ismail, Gunady, Bravo and Feizi (2020)¹⁸ reported
162 that several recurrent architectures failed to provide useful attributions for the time-series data. They further reported that some
163 architectures could extract meaningful time steps but fail to identify noteworthy features within those time steps. Results show
164 that our *whole MILC* model resolves the vanishing saliency problem and is a good tool for introspection of the multivariate
165 signal dynamics.

166 Interpretation of deep learning models may uncover domain-specific knowledge^{36,37} that would otherwise require high cost,
167 effort, and time investments. Often, it may also assist in identifying if the model has inherited any inherent bias from the data.
168 On the other hand, some studies^{38,39} raised doubts about the transparency of deep learning models and the applicability of
169 popular interpretability methods. Notwithstanding these diverging opinions, the significance of interpretability and visualization

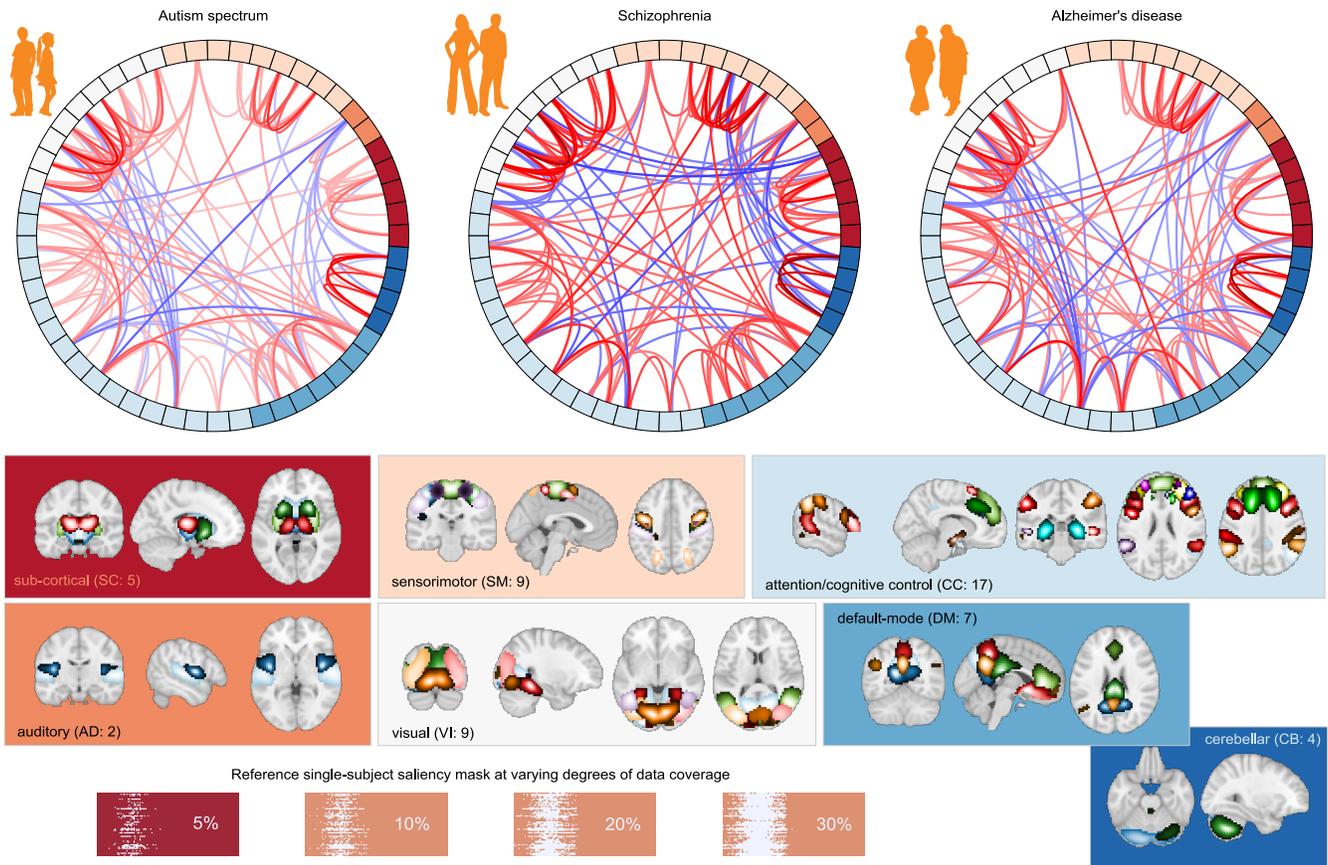


Figure 4. Top 10% FNC for patients computed using most 5% of the salient data as thresholded using feature attribution maps (saliency maps) for different disorders. Apart from the high predictive capacity of the salient data, we observed some intriguing differences among these connectograms. The autism spectrum disorder exhibits the lowest between-domain FNC. However, salient data in autism disorder highlights domain changes in specific cerebellum, sensorimotor, and subcortical domains. The model-identified salient data reflects the most widespread pattern for schizophrenia and is consistent with the literature showing cerebellum interaction across multiple domains and sensorimotor changes. The predictive features for Alzheimer's disease disease mainly concentrate on visual and cognitive interactions.

170 in medicine and healthcare cannot be overstated⁴⁰ and should involve medical experts as well. Expert human involvement in
 171 interpreting the extracted information on clinical terms may help validate and guide disease-associated discovery. A recent
 172 review⁴¹ reveals that deep learning models are a viable clinical supportive tool in the neuroimaging domain. However, studies
 173 have concentrated mainly on structural imaging data. Conversely, this paper introspects deep learning models for multivariate
 174 time-series data, which we think is an essential step toward interpretability research of functional imaging data. To this end, our
 175 model introspection results reveal the capacity of the proposed model to locate highly predictive disease-relevant information.
 176 Specifically, we validate the efficacy of the estimated feature attributions by proposing a method called **RAR**. With **RAR** and
 177 an independent SML model, we verify that IG and SGIG, when applied to *whole MILC* model, are robust, stable, and can
 178 demonstrably identify disorder-relevant parts of the brain dynamics. Precisely, the model-identified features offer very high
 179 predictive performance compared to random baselines for schizophrenia, Alzheimer's disease, and autism spectrum disorders.
 180 Moreover, our FNC analysis on model introspection results, as shown in Figure 5, harmonizes with the prior work³²⁻³⁴ for all
 181 the disorders.

182 We analyzed the required "what" and "when" aspects of the discriminative dynamics the model captured for patient-control
 183 distinction. Toward this goal, FNC analysis on the salient data revealed the minimally required connectivity ("what") of
 184 the discriminative dynamics that the model used to distinguish patients from controls. We further investigated if the model
 185 leveraged any temporal ("when") information for its discriminating power. Accordingly, we analyzed when, if such information
 186 exists, the discriminative events happen and how this temporal behavior changes between patients and controls for each disorder.
 187 As such, we analyzed the temporal densities computed from salient 5% data. Interestingly, for schizophrenia and Alzheimer's
 188 disorders, we observed that the model used temporally dense information to distinguish patients from controls. However, no

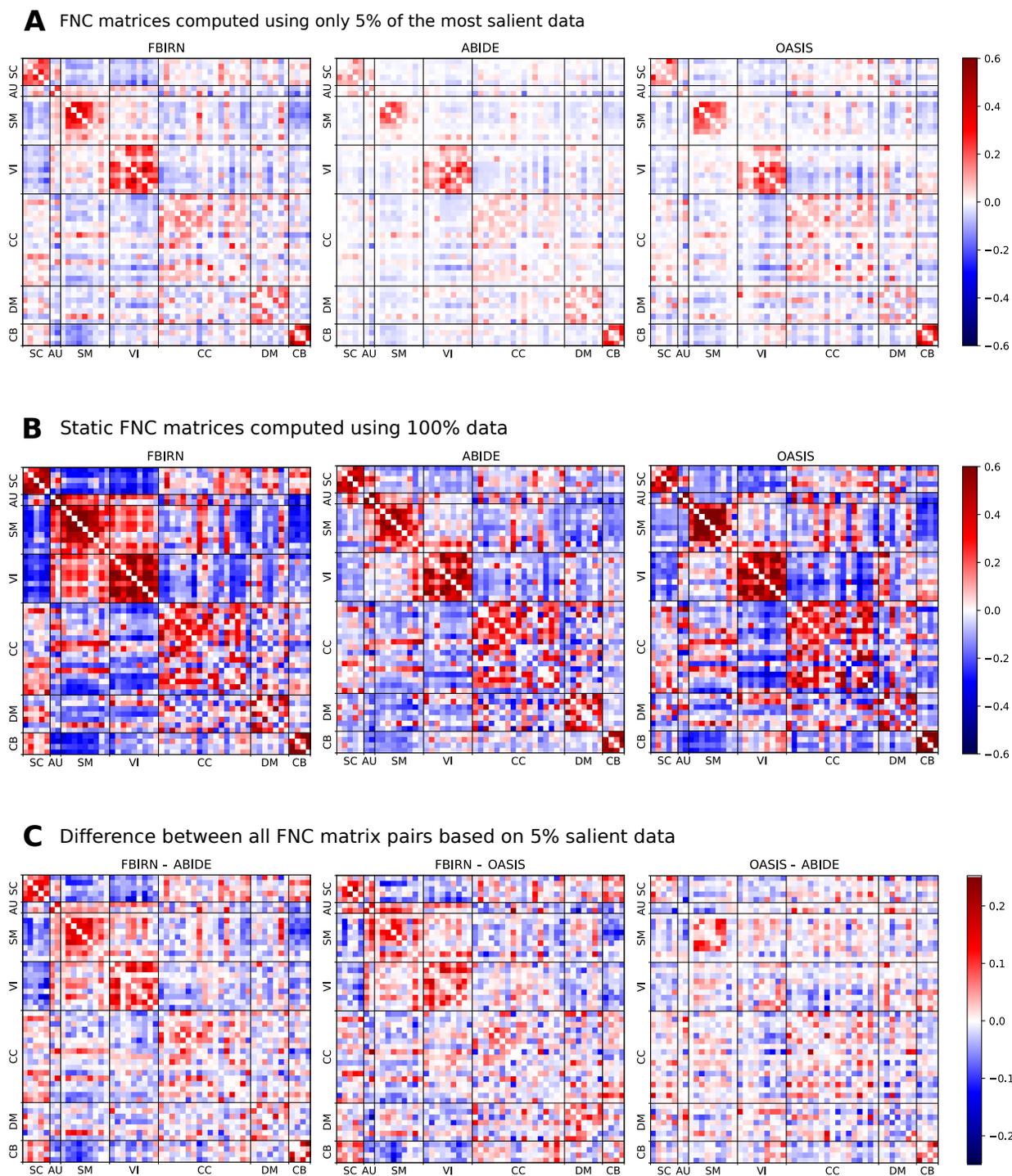
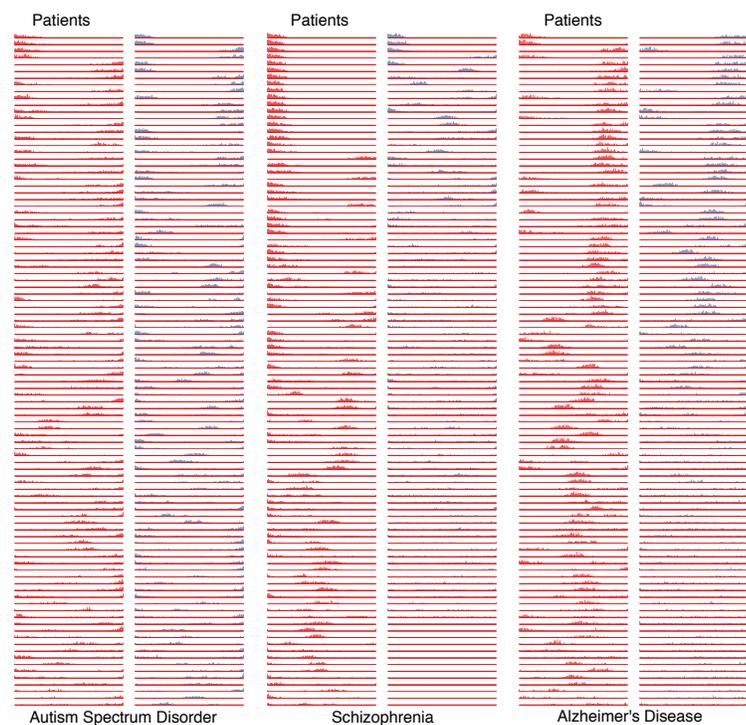


Figure 5. A: Full FNC for patients computed using most 5% of the salient data selected based on feature attribution values for different disorders. **B:** Static FNC (i.e., using 100% data) matrices for patients of different disorders. The FNC based on 5% salient data (**A**) does indeed convey the same focused dynamic information as currently assessed in FNC matrices based on 100% data (**B**). It is thus apparent that the proposed model can capture the focused information aligned with the current domain knowledge. **C:** Pairwise difference of FNC matrices based on 5% salient data. The difference FNC matrices based on focused data indicate that each disorder has a uniquely distinguishable association with brain dynamics.

A Temporal densities based on 5% salient data



B Uniformity/spikiness distributions of temporal densities

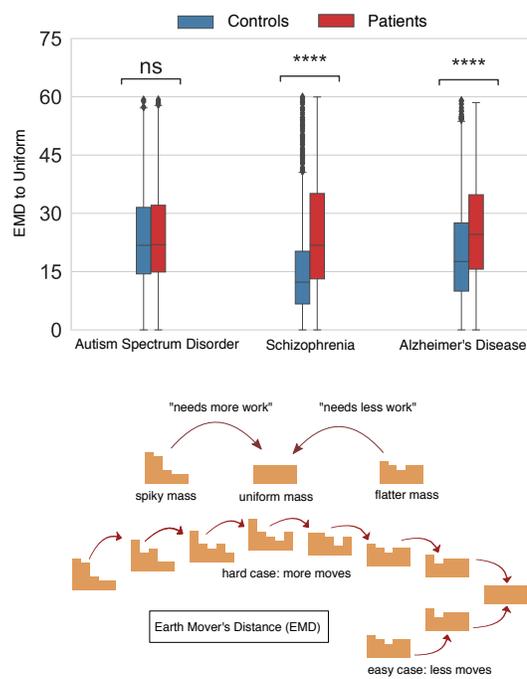


Figure 6. A: Examples of the temporal density based on the top 5% values of the saliency maps from patients and controls for each disorder. It is noticeable that the temporal density for schizophrenia and Alzheimer's patients is more focal in time as reflected in the spikiness, indicating that the discriminative activity for patients occurs predominantly in a shorter time interval. In contrast, for controls, model predictions do not relate to specific time intervals. For autism spectrum disorder, however, the *whole MILC* model did not capture any temporal adherence to the discriminative activity for patients. That is, the discriminatory events are not focal on shorter time intervals for ASD. **B:** The EMD (Earth Mover's Distance) distributions as a proxy measure for uniformity/spikiness of temporal densities. We analyzed the EMD measures of patients and controls to investigate the discriminative properties of salient data in terms of the spikiness or uniformity of the temporal densities. The larger EMD measures for schizophrenia and Alzheimer's patients substantiate that the model found the discriminative activity in shorter focused time intervals. In contrast, for ASD, the equal EMD values for both patients and controls indicate that the temporal density measures do not relate to the discriminative activity for this disorder.

189 temporal association is noticed in the model behavior to distinguish ASD patients from controls. We substantiate this aspect of
 190 temporal association using a non-parametric statistical test as shown in Figure 6.

191 Deep learning models typically require large amounts of data for efficient training. However, in the field of neuroimaging,
 192 collecting massive amounts of homogeneous data is infeasible thus constraining researchers to work with small data. In such
 193 cases, transfer learning²⁰⁻²³ is practically helpful to enable learning directly from data. Self-supervised learning has made
 194 significant progress in computer vision classification tasks²⁴ and is equally applicable to deep convolutional and recurrent
 195 networks. As demonstrated, our self-supervised pretraining scheme²² enables downstream learning with minimal training data,
 196 making the direct investigation of system dynamics feasible. Our findings demonstrate that self-supervised pretraining on
 197 healthy adults dataset noticeably uplifts the downstream model's performance on a disparate disorder dataset. These benefits
 198 generalize across datasets and disorders and thus alleviate the need to collect a massive amount of expensive data.

199 To conclude, we interpret DL models trained on fMRI signals to classify mental disorders from healthy controls to provide
 200 means to identify salient parts of the brain dynamics (activity patterns). In particular, we show that one can capture the dynamic
 201 signatures as generally captured in traditional full data functional network connectivity (FNC). We further demonstrate that
 202 the brain function manifests itself via unique dynamic signatures across time scales (latent temporality) in various disorders.
 203 Subsequently, we present an adaptive, interpretable methodology to capture these temporally transient dynamic signatures
 204 that can help distinguish disorders. Understanding the spatial and temporal specificity of the brain activity patterns will help
 205 establish the technique for clinical use by relating the differences in signature to symptoms. Moreover, to achieve these

desirable disorder-specific insights, the proposed pretraining method waives the need for well-defined ground truth (biomarkers) about the disorder under consideration and a larger sample size. In the future, this method could be a significant step towards establishing more robust correlates of function-structure dependency in the brain and can also be applied more broadly to understand inter-and intraindividual variability and alterations across psychiatric disorders.

Methods

The proposed methodology consists of 4 steps: model pretraining, downstream classification, feature importance estimation, and feature evaluation. First, we pre-trained the proposed network (*whole MILC*)²² on a large unrelated and unlabeled dataset to learn valuable latent representations. This pretraining, as described in the *whole MILC* Section, intuitively lets the network learn foundational knowledge about the dynamics only from the healthy subjects. For pretraining and downstream tasks, we used the same model as used in²². However, for the current study, we replaced the CNN encoder with a recurrent encoder because we found it more stable for post hoc explanations of multivariate time-series data while interpreting the model’s predictions. As the learned dynamics are directly transferable, we used the pre-trained network to discriminate patients from healthy controls in different downstream tasks. In the second step, we trained the downstream classification model to learn more from the downstream training data dynamics. In the third step, we estimated feature importance values based on the model’s predictions using different interpretability methods (see Model Interpretability section). In the fourth step, we evaluated the estimated features using *RAR* method and an SVM model as described in the *RAR* Section. Before going through the methodological pipeline, we preprocessed the data as described below:

Preprocessing

We preprocessed the raw resting-state fMRI data using statistical parametric mapping (SPM12, <http://www.fil.ion.ucl.ac.uk/spm/>) MATLAB package. After the preprocessing, we selected those subjects in the analysis which have head motions $\leq 3^\circ$ and ≤ 3 mm. To ensure high data quality, we performed quality control (QC) on the spatial normalization output and removed subjects with limited brain coverage⁴². For each dataset, we used ICA components derived via a fully automated approach⁴³ using the same procedure as described in⁴². We used ICA time courses as these offer a better representation of the data than anatomical or fixed atlas-based approaches⁴⁴. This study used 53 intrinsic networks (components) for all experiments. In pretraining, we used a sliding window of 53×20 size with stride = 10 along the time dimension to feed the ICA time courses through a parameter-shared encoder. In all downstream classification experiments, we used a similar sliding window with stride = 1.

Whole MILC

The *whole MILC* model, as shown in Figure 7, consists of two unidirectional LSTM models arranged in a top-down fashion. While the low-level LSTM functioned as a parameter-shared encoder for the sliding window over ICA time courses, the top-level LSTM used the encoder embeddings to generate a global representation for the entire sequence. Both LSTM models separately applied an attention mechanism⁴⁵ to retain interpretable information for further model introspection. One of the benefits of the *whole MILC* model is that it is pre-trainable. Moreover, the learned representations are directly transferable to a set of downstream discriminative tasks. The *whole MILC* model used a self-supervised pretraining objective²² that maximized the mutual information between the latent space of a window (time slice from ICA time courses) and the corresponding whole sequence (complete ICA time courses per subject).

Let $D = \{(\mathbf{u}_t^i, \mathbf{v}^j) : 1 \leq t \leq T, 1 \leq i, j \leq N\}$ be a dataset of window-sequence embedding pairs computed from ICA time courses, where subscript t refers to the t -th window, superscripts i, j each refers to a sequence number. T is the number of windows in a sequence, and N is the total number of sequences in the dataset. D can be decomposed into a set of positive pairs D^+ ($i = j$) and a set of negative pairs D^- ($i \neq j$) denoting a joint and a marginal distribution respectively for the window-sequence pairs in the latent space. With a separable function f , we used InfoNCE estimator⁴⁶ to compute a lower bound $\mathcal{I}_f(D^+)$ on the mutual information defined as:

$$\mathcal{I}(D^+) \geq \mathcal{I}_f(D^+) \triangleq \sum_{i=1}^N \sum_{t=1}^T \log \frac{\exp f((\mathbf{u}_t^i, \mathbf{v}^i))}{\sum_{k=1}^N \exp f((\mathbf{u}_t^i, \mathbf{v}^k))}, \quad (1)$$

f was defined as $f(\mathbf{u}_t, \mathbf{v}) = \phi(\mathbf{u}_t^i)^T(\mathbf{v}^j)$, where ϕ was some embedding function learnt by network parameters. f learned an embedding function such that it assigned higher values for positive pairs than for negative pairs, i.e., $f(D^+) \gg f(D^-)$. To make it precise, \mathbf{u}_t and \mathbf{v} in the Equation 1 respectively refer to window embedding \mathbf{z}_t and global sequence embedding \mathbf{c} in Figure 7. The InfoNCE loss using f as a representation model is defined as $L = -\mathcal{I}_f$.

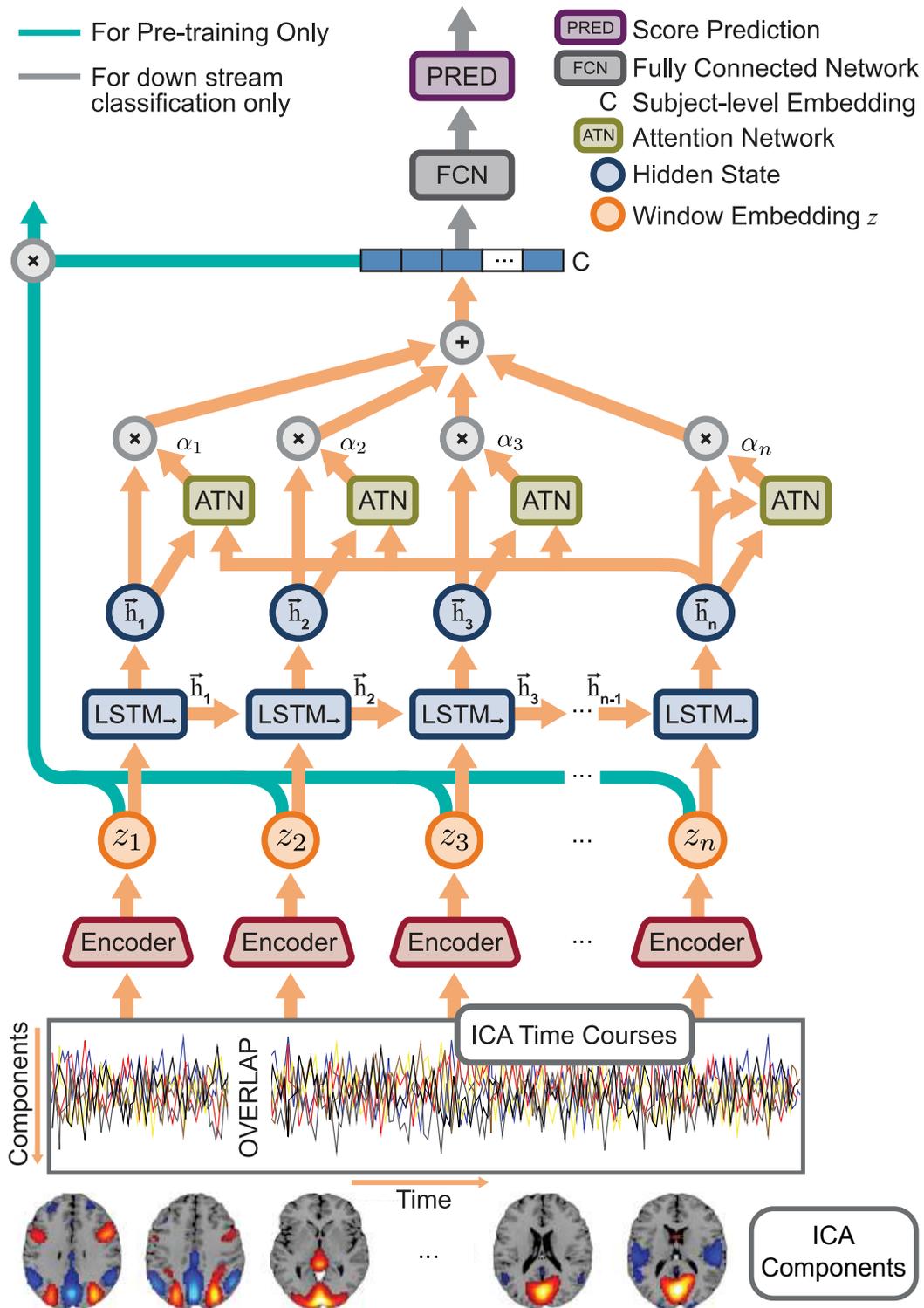


Figure 7. The *whole* MILC architecture—an attention-based top-down recurrent network. Precisely, we used an LSTM network with an attention mechanism as a parameter-shared encoder to generate the latent embeddings z for the sliding window at all relevant positions. The top LSTM network (marked as LSTM) used these embeddings (z) to obtain the global representation c for the entire subject. During pretraining, we intended to maximize the mutual information between z and c . In the downstream classification task, we used the global representation c directly as input to a fully connected network for predictions. Based on these predictions, we estimated feature attributions using different interpretability methods. Finally, we evaluated the feature attributions using the RAR method and an SVM model.

251 **Attention Mechanism**

252 The attention mechanism is a valuable construct commonly used in DL architecture to preserve long-term dependency in the
253 recurrent neural network. Initially, Bahdanau, Cho, and Bengio (2014)⁴⁵ introduced the attention mechanism for the neural
254 machine translation to compute the relevance of source words toward each output word. However, the attention mechanism
255 can benefit other applications too. For example, we used the attention mechanism to solve vanishing saliency problems in
256 the LSTM networks to retain interpretable information during model training. In the attention mechanism as used in *whole*
257 *MILC* model, we took all the hidden states $\mathbf{h} = [\mathbf{h}_1, \mathbf{h}_2, \dots, \mathbf{h}_n]$ from the LSTM network and concatenated each hidden state \mathbf{h}_i
258 with the hidden state at the last time step \mathbf{h}_n before passing through an attention mechanism f_a . The attention mechanism f_a ,
259 similar to the additive attention mechanism introduced in⁴⁵, took pairs of hidden states $(\mathbf{h}_i, \mathbf{h}_n)$ as inputs, passed through a
260 2-layer feed-forward network and generated a vector of n alignment scores $f_a(\mathbf{h}_i, \mathbf{h}_n)$. The alignment score for each time point i
261 intuitively indicates the degree of relevance of the corresponding hidden state to the overall embedding. We normalized the
262 alignment scores using softmax to produce a series of weights $\alpha_1, \alpha_2, \dots, \alpha_n$. α_i for each time point is defined as:

$$\alpha_i = \frac{\exp(f_a(\mathbf{h}_i, \mathbf{h}_n))}{\sum_{k=1}^n \exp(f_a(\mathbf{h}_k, \mathbf{h}_n))} \quad (2)$$

where n was the number of time steps over which attention was applied. Note that the value of n for the encoder LSTM network (for the sliding window) differed from the top LSTM network (for the full subject). The global representation \mathbf{c} (or the window embedding \mathbf{z}) was generated using the formula as follows:

$$\mathbf{c} = \sum_{k=1}^n \alpha_k \mathbf{h}_k \quad (3)$$

263 **whole MILC Setup**

264 **Encoder Embedding:** The LSTM encoder with an attention mechanism used a sliding window of 53×20 size to feed the ICA
265 time courses and encoded features at each time point into a 256-dimensional representation. At each position of the sliding
266 window, we concatenated the hidden state for each time step t_i within the window with the final hidden state of the same window
267 as described in the attention mechanism. We then passed these concatenated 512-dimensional vectors through an attention
268 network, a two-layer feed-forward network with hidden units 64. The network learned a series of weights representatives of the
269 magnitude of attention regarded as important for the time steps. All the hidden representations within a window were then
270 weighted based on the attention scales to produce window embedding \mathbf{z} .

271 **Pretraining:** In *whole MILC* based pretraining, we passed all the encoder embeddings $\mathbf{z} = \mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_n$ to another unidirectional
272 LSTM network with an attention mechanism. In this top recurrent network, each window embedding \mathbf{z}_i corresponded to
273 the input for a single time step. We used 200 dimensions to represent the hidden state for this top network. We concatenated
274 each hidden state with the hidden state at the last time step to make it contextually relevant for the attention mechanism. The top
275 attention network used 400 input neurons and 128 hidden units to learn k weights, where k was the number of input windows.
276 These weights were used as coefficients in the linear combination of hidden representations to generate a global embedding \mathbf{c}
277 of dimension 200 for each subject. Based on \mathbf{c} and \mathbf{z} , we pre-trained the neural network to maximize the mutual information
278 between a window and the corresponding input sequence. We used subjects from the HCP dataset for pretraining and used 700
279 subjects for training and 123 subjects for the test, obtaining 89% pretraining accuracy.

280 **Classification Tasks:** In downstream tasks, we deal with classifying subjects into patients and controls separately for each
281 disorder. Similar to pretraining, we fed ICA time courses into the LSTM encoder using a sliding window. The LSTM encoder
282 projected all the windows into latent representations \mathbf{z} , which were then passed to another LSTM network to obtain a global
283 representation \mathbf{c} . Finally, on top of \mathbf{c} , we used a feed-forward network with 200 hidden units to perform binary classification.
284 We gradually increased the number of supervised training subjects to observe the pretraining effect on downstream data size
285 compared to the setup where we used no pretraining. For each experiment, we report cross-validated results. Moreover, we
286 performed ten repetitions of each experimental setup, with different random seeds for every cross-validation fold to ensure
287 stable results. For each random seed, we randomly chose the training samples as required from the available training pool.

288 **Model Interpretability**

289 The need to enable model interpretation led to a variety of model introspection techniques that can be roughly split into three
290 groups: 1) model-sensitive^{25,26}, 2) model-agnostic^{47,48}, and 3) counterfactual explanations⁴⁹. The techniques have their
291 relative benefits and pitfalls in addressing the desiderata of different applications⁵⁰. Adebayo, Muelly, Liccardi, and Kim
292 (2020)⁵¹ reported that, under normal conditions, gradients, smoothgrad²⁶, and integrated gradients (IG)²⁵ passed end-user
293 recommendations. Additionally, the smoothgrad method²⁶ resolves the problems⁵² of saliency maps, which in general, are
294 susceptible to noise and input perturbations. Guided by these findings, we relied on IG, and smoothgrad on IG to introspect the
295 proposed model. Notably, we found IG and smoothgrad on IG generalizable, stable, and noise-robust across the disorders.

296 **Random Baseline**

297 We randomly assigned feature importance values to create random baselines to validate the post hoc explanations (saliency
 298 maps). Specifically, we ordered the features uniformly at random using random permutations and considered each permutation
 299 as an order of importance. We refer to this random estimator as g^R throughout the paper. In contrast, we used the magnitude
 300 of the estimated attribution values as the order of importance for the model-generated post hoc explanations. To evaluate the
 301 efficacy of the estimated feature importance, we compared the predictive power of the model-estimated salient features against
 302 random baselines using a technique called **RAR**, which we describe below.

303 **RAR Method and Setup**

304 In **RAR**, we retained only a small percentage of the most salient features as determined by the model and replaced other
 305 features with non-informative values (zeros). We used these modified samples to retrain an SVM model to evaluate the
 306 effectiveness of the estimated feature attributions. In particular, we show that the performance obtained with *whole MILC*
 307 model-estimated salient features far exceeded the random baseline. We mathematically describe the **RAR** scheme as follows:

308 Let us define X to be the original dataset. $X^M | g^R$ be the modified dataset based on random importance estimates and
 309 $X^M | g_i$ be the modified dataset according to the saliency maps generated by applying some interpretability method g_i on *whole*
 310 *MILC* predictions. We computed static functional network connectivity, measured as Pearson’s correlation coefficients, for
 311 each sample in X^M . We used these correlation coefficients as features to train an independent SVM model de novo. We
 312 evaluated the classification performance of the SVM models trained separately with *whole MILC*-generated salient features and
 313 randomly selected features. Indeed, we show that $\xi(X^M | g_i) > \xi(X^M | g^R)$, where ξ is the performance evaluation function,
 314 e.g. area under the ROC curve and/or accuracy.

315 It is to note that we sorted the features based on their signed attribution values before considering them for validation. We
 316 searched for the SVM (nonlinear) parameters using a parameter grid and 3-fold cross-validation on the training data. We used
 317 the same folds and train-test splits for the **RAR** evaluation as used in the *whole MILC* model. Figure 8 shows the schematic of
 318 the end-to-end process: 1) training the *whole MILC* and feature attributions and 2) Evaluation of the feature attributions using
 319 **RAR** and an SVM model.

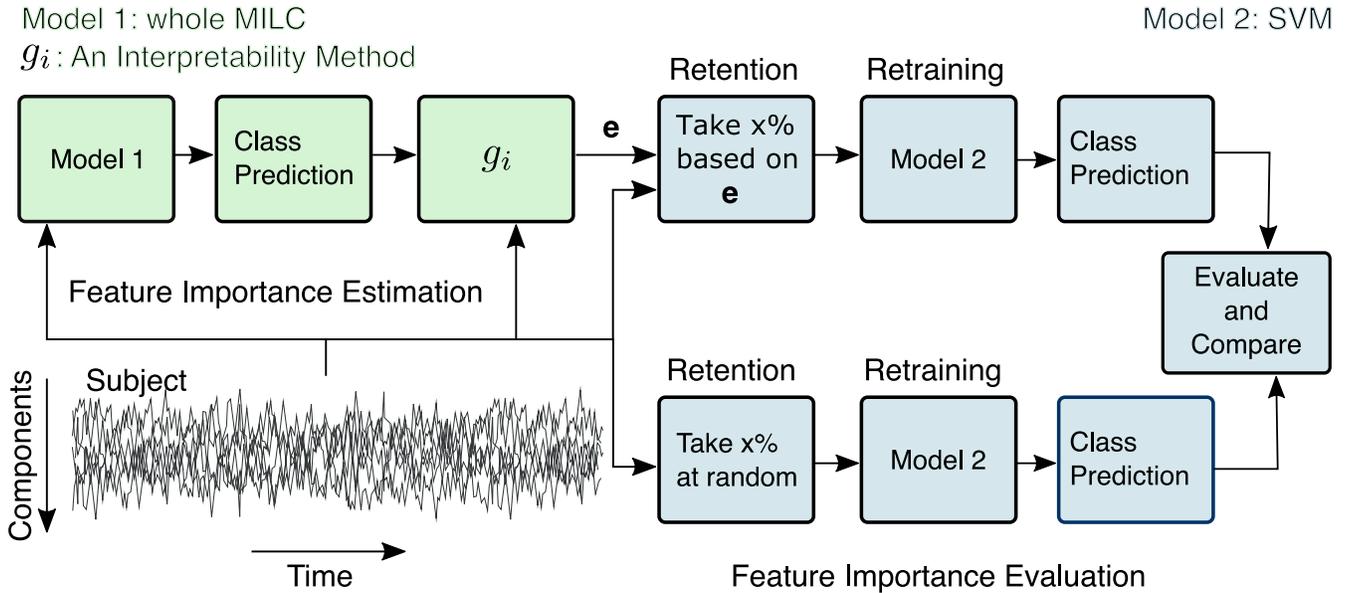


Figure 8. End-to-end process of **RAR** evaluation. For each subject in the dataset, based on the *whole MILC* class prediction and model parameters, we estimated the feature importance vector e using some interpretability method g_i . Later on, we validated these estimates against random feature attributions g^R using the **RAR** method and an SVM model. Through the SVM model’s performance when separately trained with different feature sets, we show that *whole MILC* model-estimated features were highly predictive compared to a random selection of a similar amount of features. Empirically, we show that $\xi(X^M | g_i) > \xi(X^M | g^R)$, where ξ is the performance evaluation function (e.g., area under the curve) and X^M refers to the modified dataset constructed based on only retained feature values.

References

- 320 **1.** Goldberg, D. P. & Huxley, P. *Common mental disorders: a bio-social model*. (Tavistock/Routledge, 1992).
- 321
- 322 **2.** Calhoun, V. D., Miller, R., Pearlson, G. & Adali, T. The chronnectome: time-varying connectivity networks as the next
323 frontier in fMRI data discovery. *Neuron* **84**, 262–274 (2014).
- 324 **3.** Sui, J., Jiang, R., Bustillo, J. & Calhoun, V. Neuroimaging-based individualized prediction of cognition and behavior for
325 mental disorders and health: methods and promises. *Biol. psychiatry* (2020).
- 326 **4.** Hutchison, R. M. *et al.* Dynamic functional connectivity: promise, issues, and interpretations. *Neuroimage* **80**, 360–378
327 (2013).
- 328 **5.** Logothetis, N. K. What we can do and what we cannot do with fMRI. *Nature* **453**, 869–878 (2008).
- 329 **6.** Heeger, D. J. & Ress, D. What does fMRI tell us about neuronal activity? *Nat. Rev. Neurosci.* **3**, 142–151 (2002).
- 330 **7.** Bellman, R. Adaptive control processes: a guided tour princeton university press. *Princeton, New Jersey, USA* 96 (1961).
- 331 **8.** Khazaei, A., Ebrahimzadeh, A. & Babajani-Feremi, A. Application of advanced machine learning methods on resting-state
332 fMRI network for identification of mild cognitive impairment and Alzheimer’s disease. *Brain Imaging Behav.* **10**, 799–817,
333 DOI: [10.1007/s11682-015-9448-7](https://doi.org/10.1007/s11682-015-9448-7) (2016).
- 334 **9.** Straathof, M., Sinke, M. R., Dijkhuizen, R. M. & Otte, W. M. A systematic review on the quantitative relationship between
335 structural and functional network connectivity strength in mammalian brains. *J. Cereb. Blood Flow & Metab.* **39**, 189–209
336 (2019).
- 337 **10.** Hummer, T. A. *et al.* Functional network connectivity in early-stage schizophrenia. *Schizophr. research* **218**, 107–115
338 (2020).
- 339 **11.** Liang, X. *et al.* Effects of different correlation metrics and preprocessing factors on small-world brain functional networks:
340 a resting-state functional mri study. *PloS one* **7**, e32766 (2012).
- 341 **12.** Zhang, Y., Zhang, H., Chen, X., Lee, S.-W. & Shen, D. Hybrid high-order functional connectivity networks using
342 resting-state functional mri for mild cognitive impairment diagnosis. *Sci. reports* **7**, 1–15 (2017).
- 343 **13.** Plis, S. M. *et al.* Deep learning for neuroimaging: a validation study. *Front. neuroscience* **8**, 229 (2014).
- 344 **14.** Vieira, S., Pinaya, W. H. & Mechelli, A. Using deep learning to investigate the neuroimaging correlates of psychiatric and
345 neurological disorders: Methods and applications. *Neurosci. & Biobehav. Rev.* **74**, 58–75 (2017).
- 346 **15.** Zhang, L., Wang, M., Liu, M. & Zhang, D. A survey on deep learning for neuroimaging-based brain disorder analysis.
347 *Front. neuroscience* **14** (2020).
- 348 **16.** Lewis, N. *et al.* (in press) can recurrent models know more than we do? In *2021 IEEE International Conference on*
349 *Healthcare Informatics (ICHI)* (2021).
- 350 **17.** Abrol, A. *et al.* Deep learning encodes robust discriminative neuroimaging representations to outperform standard machine
351 learning. *Nat. communications* **12**, 1–17 (2021).
- 352 **18.** Ismail, A., Gunady, M., Bravo, H. & Feizi, S. Benchmarking deep learning interpretability in time series predictions. *Adv.*
353 *Neural Inf. Process. Syst. Foundation (NeurIPS)* (2020).
- 354 **19.** Sun, C., Shrivastava, A., Singh, S. & Gupta, A. Revisiting unreasonable effectiveness of data in deep learning era. In
355 *Proceedings of the IEEE international conference on computer vision*, 843–852 (2017).
- 356 **20.** Mensch, A., Mairal, J., Bzdok, D., Thirion, B. & Varoquaux, G. Learning neural representations of human cognition across
357 many fMRI studies. In *Advances in Neural Information Processing Systems*, 5883–5893 (2017).
- 358 **21.** Thomas, A. W., Müller, K.-R. & Samek, W. Deep transfer learning for whole-brain fMRI analyses. *arXiv preprint*
359 *arXiv:1907.01953* (2019).

- 360 **22.** Mahmood, U. *et al.* Whole MILC: generalizing learned dynamics across tasks, datasets, and populations. In *Proceedings*
361 *of the 23rd international conference on medical image computing and computer assisted intervention (MICCAI)*, 407–417
362 (Springer, 2020).
- 363 **23.** Mahmood, U., Rahman, M. M., Fedorov, A., Fu, Z. & Plis, S. Transfer learning of fMRI dynamics. *arXiv preprint*
364 *arXiv:1911.06813* (2019).
- 365 **24.** Newell, A. & Deng, J. How useful is self-supervised pretraining for visual tasks? In *Proceedings of the IEEE/CVF*
366 *Conference on Computer Vision and Pattern Recognition*, 7345–7354 (2020).
- 367 **25.** Sundararajan, M., Taly, A. & Yan, Q. Axiomatic attribution for deep networks. *arXiv preprint arXiv:1703.01365* (2017).
- 368 **26.** Smilkov, D., Thorat, N., Kim, B., Viégas, F. & Wattenberg, M. Smoothgrad: removing noise by adding noise. *arXiv*
369 *preprint arXiv:1706.03825* (2017).
- 370 **27.** Van Essen, D. C. *et al.* The WU-Minn human connectome project: an overview. *Neuroimage* **80**, 62–79 (2013).
- 371 **28.** Di Martino, A. *et al.* The autism brain imaging data exchange: towards a large-scale evaluation of the intrinsic brain
372 architecture in autism. *Mol. psychiatry* **19**, 659 (2014).
- 373 **29.** Keator, D. B. *et al.* The function biomedical informatics research network data repository. *Neuroimage* **124**, 1074–1079
374 (2016).
- 375 **30.** Rubin, E. H. *et al.* A prospective study of cognitive function and onset of dementia in cognitively healthy elders. *Arch.*
376 *neurology* **55**, 395–401 (1998).
- 377 **31.** Calhoun, V. D., Adali, T., Pearlson, G. D. & Pekar, J. A method for making group inferences from functional MRI data
378 using independent component analysis. *Hum. brain mapping* **14**, 140–151 (2001).
- 379 **32.** Monk, C. S. *et al.* Abnormalities of intrinsic functional connectivity in autism spectrum disorders. *Neuroimage* **47**,
380 764–772 (2009).
- 381 **33.** Li, S. *et al.* Dysconnectivity of multiple brain networks in schizophrenia: a meta-analysis of resting-state functional
382 connectivity. *Front. psychiatry* **10**, 482 (2019).
- 383 **34.** Gu, Y. *et al.* Abnormal dynamic functional connectivity in alzheimer’s disease. *CNS neuroscience & therapeutics* **26**,
384 962–971 (2020).
- 385 **35.** Rubner, Y., Tomasi, C. & Guibas, L. J. The earth mover’s distance as a metric for image retrieval. *Int. journal computer*
386 *vision* **40**, 99–121 (2000).
- 387 **36.** Hicks, S. A. *et al.* Explaining deep neural networks for knowledge discovery in electrocardiogram analysis. *Sci. Reports*
388 **11**, 1–11 (2021).
- 389 **37.** Ghorbani, A. *et al.* Deep learning interpretation of echocardiograms. *NPJ digital medicine* **3**, 1–10 (2020).
- 390 **38.** Kindermans, P.-J. *et al.* The (un) reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing*
391 *Deep Learning*, 267–280 (Springer, 2019).
- 392 **39.** Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models
393 instead. *Nat. Mach. Intell.* **1**, 206–215 (2019).
- 394 **40.** Vellido, A. The importance of interpretability and visualization in machine learning for applications in medicine and health
395 care. *Neural computing applications* 1–15 (2019).
- 396 **41.** Singh, A., Sengupta, S. & Lakshminarayanan, V. Explainable deep learning models in medical image analysis. *J. Imaging*
397 **6**, 52 (2020).
- 398 **42.** Fu, Z. *et al.* Altered static and dynamic functional network connectivity in alzheimer’s disease and subcortical ischemic
399 vascular disease: shared and specific brain connectivity abnormalities. *Hum. Brain Mapp.* (2019).
- 400 **43.** Du, Y. *et al.* NeuroMark: An automated and adaptive ICA based pipeline to identify reproducible fMRI markers of brain
401 disorders. *NeuroImage: Clin.* **28**, 102375 (2020).

- 402 **44.** Yu, Q. *et al.* Comparing brain graphs in which nodes are regions of interest or independent components: A simulation
403 study. *J. neuroscience methods* **291**, 61–68 (2017).
- 404 **45.** Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly learning to align and translate. *arXiv preprint*
405 *arXiv:1409.0473* (2014).
- 406 **46.** Oord, A. v. d., Li, Y. & Vinyals, O. Representation learning with contrastive predictive coding. *arXiv preprint*
407 *arXiv:1807.03748* (2018).
- 408 **47.** Castro, J., Gómez, D. & Tejada, J. Polynomial calculation of the shapley value based on sampling. *Comput. & Oper. Res.*
409 **36**, 1726–1730 (2009).
- 410 **48.** Ribeiro, M. T., Singh, S. & Guestrin, C. "Why should i trust you?" explaining the predictions of any classifier. In
411 *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, 1135–1144
412 (2016).
- 413 **49.** Mothilal, R. K., Sharma, A. & Tan, C. Explaining machine learning classifiers through diverse counterfactual explanations.
414 In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 607–617 (2020).
- 415 **50.** Samek, W., Montavon, G., Lapuschkin, S., Anders, C. J. & Müller, K.-R. Toward interpretable machine learning:
416 Transparent deep neural networks and beyond. *arXiv preprint arXiv:2003.07631* (2020).
- 417 **51.** Adebayo, J., Muelly, M., Liccardi, I. & Kim, B. Debugging tests for model explanations. In Larochelle, H., Ranzato, M.,
418 Hadsell, R., Balcan, M. F. & Lin, H. (eds.) *Advances in Neural Information Processing Systems*, vol. 33, 700–712 (Curran
419 Associates, Inc., 2020).
- 420 **52.** Levine, A., Singla, S. & Feizi, S. Certifiably robust interpretation in deep learning. *arXiv preprint arXiv:1905.12105*
421 (2019).

422 **Acknowledgements**

423 This study was supported by startup funds to SMP and in part by NIH grants R01EB006841, R01MH118695, and RF1MH121885.
424 Data for healthy subjects was provided [in part] by the Human Connectome Project, WU-Minn Consortium (Principal In-
425 vestigators: David Van Essen and Kamil Ugurbil; 1U54MH091657) funded by the 16 NIH Institutes and Centers that
426 support the NIH Blueprint for Neuroscience Research; and by the McDonnell Center for Systems Neuroscience at Wash-
427 ington University. Data for Schizophrenia used in this study were downloaded from the Function BIRN Data Reposi-
428 tory (<http://bdr.birncommunity.org:8080/BDR/>), supported by grants to the Function BIRN (U24-RR021992)
429 Testbed funded by the National Center for Research Resources at the National Institutes of Health, U.S.A. Data for Alzheimer's
430 was provided by OASIS-3: Principal Investigators: T. Benzinger, D. Marcus, J. Morris; NIH P50AG00561, P30NS09857781,
431 P01AG026276, P01AG003991, R01AG043434, UL1TR000448, R01EB009352. AV-45 doses were provided by Avid Radio-
432 pharmaceuticals, a wholly owned subsidiary of Eli Lilly. Autism data was provided by ABIDE. We acknowledge primary
433 support for the work by Adriana Di Martino provided by the (NIMH K23MH087770) and the Leon Levy Foundation and
434 primary support for the work by Michael P. Milham and the INDI team was provided by gifts from Joseph P. Healy and the
435 Stavros Niarchos Foundation to the Child Mind Institute, as well as by an NIMH award to MPM (NIMH R03MH096321).

436 **Author contributions statement**

437 Conceptualization: M.M.R., U.M., A.F, N.L, and S.M.P. Funding acquisition: S.M.P. Supervision: S.M.P. Data collection
438 and preprocessing: Z.F. Methodology: M.M.R, U.M, A.F, N.L, and S.M.P. Software: M.M.R. Validation: M.M.R, S.M.P.
439 Visualization: M.M.R, S.M.P. Writing—original draft: M.M.R, S.M.P. Writing Revisions: M.M.R, S.M.P, V.C, and H.G. All
440 authors reviewed the manuscript.

441 **Competing Interests statement**

442 The authors do not have any competing interests.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryInformationforSubmission.pdf](#)