Check for updates

1      **A Hybrid Spatio-Temporal Deep Belief Network and Sparse Representation-Based**

2      **Framework Reveals Multi-Level Core Functional Components in Decoding Multi-Task**

3      **fMRI Signals**

4      *Limei Song[1#], Yudan Ren[1#*], Shuhan Xu[1], Yuqing Hou[1], Xiaowei He[1]*

5      [1] School of Information Science & Technology, Northwest University, China;

6      [#] These authors contributed equally to this work and should be considered co-first authors.

7      [*] Corresponding authors.

8      **Abstract**

9      Decoding human brain activity on various task-based functional brain imaging data is of great

10     significance for uncovering the functioning mechanism of the human mind. Currently, most

11     feature extraction model-based methods for brain state decoding are shallow machine learning

12     models, which may struggle to capture complex and precise spatio-temporal patterns of brain

13     activity from the highly noisy fMRI raw data. Moreover, although decoding models based on

14     deep learning methods benefit from their multi-layer structure that could extract spatio-

15     temporal features at multi-scale, the relatively large populations of fMRI datasets are

16     indispensable and the explainability of their results is elusive. To address the above problems,

17     we proposed a computational framework based on hybrid spatio-temporal deep belief network

18     and sparse representations to differentiate multi-task fMRI (tfMRI) signals. Using a relatively

19     small cohort of tfMRI data as a testbed, our framework can achieve an average classification

20     accuracy of 97.86% and define the multi-level temporal and spatial patterns of multiple

21     cognitive tasks. Intriguingly, our model can characterize the key components for differentiating

22    the multi-task fMRI signals. Overall, the proposed framework can identify the interpretable

23    and discriminative fMRI composition patterns at multiple scales, offering an effective

24    methodology for basic neuroscience and clinical research with relatively small cohorts.

25    **Keywords**: Multi-task classification, Task-based fMRI, Deep belief network, Sparse

26    representation, Functional brain network.

27    **Introduction**

28    For years, researchers have been attempting to decode the human brain states based on

29    functional magnetic resonance imaging (fMRI) data (Haynes & Rees, 2006; Jang, Plis, Calhoun,

30    & Lee, 2017; Rubin et al., 2017; Stanislas Dehaene, 1998), where distinguishing different

31    cognitive tasks from fMRI data and extracting discriminative fMRI composition patterns are

32    effective means to improve our understanding of the relationship among current cognitive tasks,

33    brain responses, and individual behavior (Friston, 2009; Logothetis, 2008). To decode

34    meaningful neurological patterns embedded in diverse task-based fMRI data, various

35    computational and statistical methods have been proposed in the last decades. The most widely

36    used brain state decoding strategy is multi-voxel pattern analysis (MVPA) (Davatzikos et al.,

37    2005; Jang et al., 2017; Kriegeskorte & Bandettini, 2007). Despite its popularity, its commonly-

38    used classification strategy support vector machine (SVM) usually struggles to perform well

39    on high-dimensional fMRI data and thus requires effective techniques for feature

40    selection/extraction (LeCun, Bengio, & Hinton, 2015; Vieira, Pinaya, & Mechelli, 2017).

41    Hence, the feasibility of feature selection/extraction has been investigated using various

42    machine learning methods (LeCun et al., 2015; Vieira et al., 2017; S. Zhang et al., 2016).

2

43    However, most of these machine learning methods rely on shallow models, and their shallow

44    nature may hinder them from effectively capturing non-linear relationships in the highly noisy

45    fMRI raw data, resulting in difficulties in extracting complex and specific spatio-temporal

46    features (Qiang et al., 2020; Rashid, Singh, & Goyal, 2020; Varoquaux & Thirion, 2014).

47        Recently, studies applying deep learning models such as deep neural network (DNN) and

48    convolutional neural networks (CNN) to decode brain states based on task-based fMRI signals

49    have been reported (J. Hu et al., 2019; Liu, He, Chen, & Gao, 2019; Sotetsu Koyamadaa, 2015;

50    Y. Zhang, Tetrel, Thirion, & Bellec, 2021). Such deep learning models take the advantage of

51    being a multi-layer architecture by stacking multiple building blocks with similar structure,

52    which has demonstrated the ability to significantly reduce noises in raw fMRI data and model

53    the non-linear relationships among neural activities of brain regions, allowing for the extraction

54    of multi-level spatio-temporal features (Bengio, Courville, & Vincent, 2012; Najafabadi et al.,

55    2015; Ren, Xu, Tao, Song, & He, 2021). Nevertheless, there are still some limitations in current

56    brain state decoding strategies based on deep learning models. First, as large-size samples are

57    indispensable for the deep learning model, current decoding models are not suitable for small

58    datasets (Bo Liu, 2017; Litjens et al., 2017; Wang et al., 2020; Wen et al., 2018). For example,

59    Wang et al. (2020) proposed a DNN-based model for tfMRI signal classification, which

60    requires 1034 subjects, making it less practical for clinical populations. Second, most of the

61    decoding models based on deep learning are end-to-end learning and the explainability of such

62    models is elusive (J. Hu et al., 2019; LeCun et al., 2015; Wang et al., 2020). Recently, some

63    researchers have attempted to define the key components for decoding brain states using the

64    machine learning method. For example, our previous study based on sparse dictionary learning

65  has determined that the key components for multi-task classification tend to be functional brain

66  networks (FBNs) (Song, Ren, Hou, He, & Liu, 2022). Another research has shown that artifact

67  components such as movement-related artifacts are significantly more informative with respect

68  to the classification accuracy of the multi-task electroencephalogram (EEG) signals

69  (McDermott et al., 2021). However, uncovering the interpretable key features in decoding

70  tfMRI signals has received much less attention.

71      Due to the pitfalls in existing research, it is desirable to develop an appropriate framework

72  capable of identifying the interpretable and discriminative fMRI composition patterns

73  embedded in multi-task fMRI data. Thus, in this study, we aim to extract both multi-level

74  group-wise temporal features and spatial features from tfMRI signals, and define interpretable

75  classification features for multi-task fMRI data simultaneously. Recent studies have revealed

76  that the deep belief network (DBN) can effectively identify multi-layer spatial and temporal

77  features from fMRI signals (Dong, 2020; Ren et al., 2021), which is typically stacked by

78  multiple Boltzmann machine (RBM) (Geoffrey E Hinton & Sejnowski, 1986) and thus can

79  naturally act as a multi-level feature extractor. Furthermore, these prior studies have integrated

80  the least absolute shrinkage and selection operator (LASSO) regression with the DBN model,

81  indicating the efficacy of LASSO regression in extracting relevant spatial patterns. Thus, we

82  here proposed a novel two-stage feature extraction framework based on hybrid DBN and sparse

83  representations framework (DBN-SR) to decode multi-task fMRI signals with the capability of

84  extracting multi-scale deep features. Specifically, the DBN model was utilized to capture multi-

85  level group-wise temporal features, based on which the individual spatial features were

86  estimated by LASSO regression. Subsequently, a sparse representation method that combines

4

87    dictionary learning and LASSO regression was utilized to further characterize the group-wise

88    spatial features and individual spatio-temporal features for the purpose of classification. Based

89    on the correspondence between the individual classification features and the group-wise spatial

90    features, a relationship between the decoding capability of classification features and their

91    spatial patterns can be effectively established, which can facilitate the interpretation of neural

92    implications associated with the classification features. Finally, due to its strong generalization

93    capabilities in small sample sizes, SVM was employed for the multi-class classification task.

94        Our results demonstrated that the proposed framework could successfully classify seven

95    task fMRI signals on a relatively small dataset. Moreover, by taking advantage of DBN in

96    extracting mid-level and high-level features and sparse coding in brain functional network

97    representation (Lv, Jiang, Li, Zhu, Chen, et al., 2015; Ren et al., 2021; Song et al., 2022), our

98    framework could effectively characterize the multi-level spatiotemporal features embedded in

99    multi-task fMRI signals, which provides the bases to identify the interpretable key components

100   for well characterizing and differentiating multi-task signals. Overall, the proposed model can

101   disclose the underlying neural implications of key components with greater classification

102   capacity, offering an effective and interpretable methodology for decoding fMRI data.
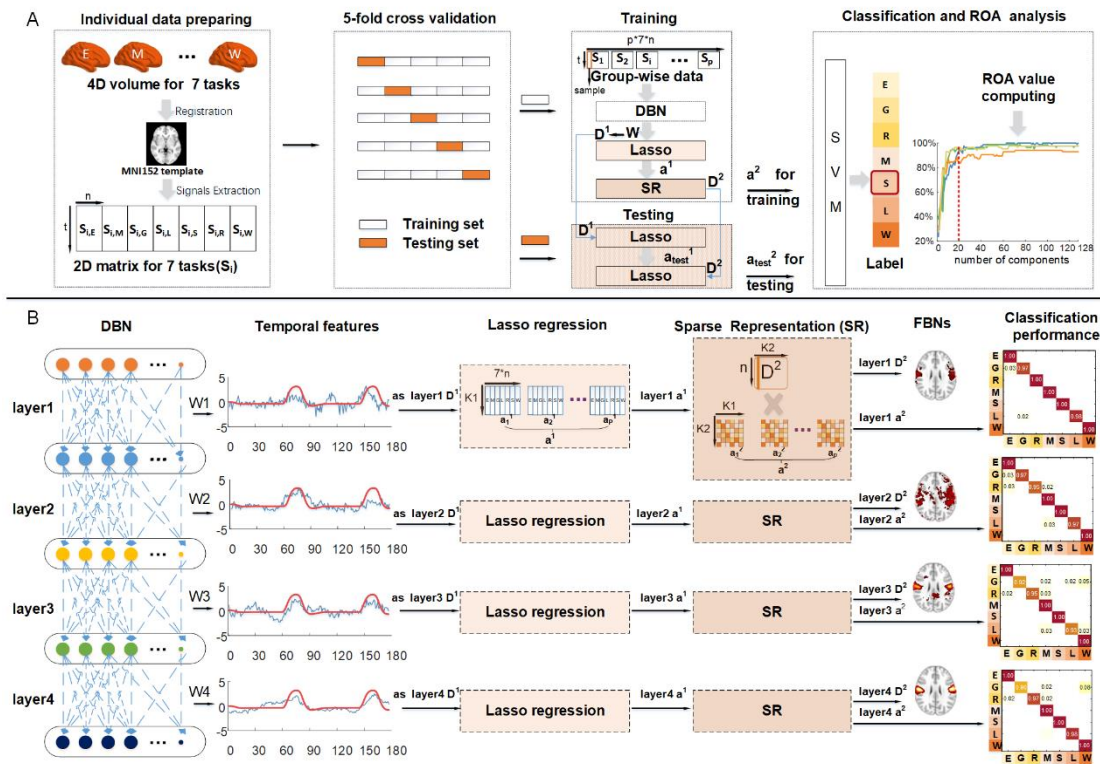
103   **Materials and methods**

104   **Overview**

105   The framework of our proposed method is illustrated in Figure 1. The pipeline of the proposed

106   framework can divide into four stages: 1) individual data preparation; 2) data preparation for

107      five-fold cross-validation; 3) training and testing process; 4) SVM-based classification and

108      Ratio of activation (ROA) analysis (Fig. 1A). In the data preparation stage, each individual's

109      tfMRI data of seven different tasks were extracted and then spatially concatenated to one signal

110      matrix (the first panel in Fig. 1A). In this work, five-fold cross-validation was performed for

111      model validation, thus the whole dataset was randomly divided into five folds (the second panel

112      in Fig. 1A). In training process, four folds were served as training set, and the tfMRI signal

113      matrices of all the subjects in training set were spatially concatenated to a multi-subject signal

114      matrix. Then, the DBN model was applied to training set to derive the weight matrix W, which

115      served as group-wise temporal features $\boldsymbol{D}^1$. Then, the LASSO regression aims to extract the

116      corresponding loading coefficient $\boldsymbol{\alpha}^1$ based on the defined temporal dictionary $\boldsymbol{D}^1$. In the

117      second stage of our model, the loading coefficient $\boldsymbol{\alpha}^1$ was employed as input to sparse

118      representations (SR) model, where they were decomposed into group-wise dictionaries $\boldsymbol{D}^2$ and

119      loading coefficient $\boldsymbol{\alpha}^2$. In testing process, the individual signal matrix in testing set and the

120      group-wise dictionary $\boldsymbol{D}^1$ obtained during the training phase was utilized as the inputs to the

121      LASSO regression. This yielded the loading coefficients $\boldsymbol{\alpha}^1_{test}$. Subsequently, employing $\boldsymbol{\alpha}^1_{test}$

122      and the $\boldsymbol{D}^2$ obtained during the training phase, we performed a second LASSO regression to

123      obtain $\boldsymbol{\alpha}^2_{test}$, which were then used as the classification features for the testing subjects (the

124      third panel in Fig. 1A). Note that during the training phase, we utilized the independent training

125      data to learn and train regularization parameters employed for LASSO regression, as well as

126      the group-wise dictionaries $\boldsymbol{D}^1$ and $\boldsymbol{D}^2$, without using any information from the test data.

127      Afterward, to further assess the multi-task fMRI data classification performance of proposed

128      model, the loading coefficient $\boldsymbol{\alpha}^2$ derived from training set was used to train support vector

6

129 machine (SVM) for classification, where the loading coefficient $\boldsymbol{\alpha}_{test}^2$ derived from testing set

130 was then fed into this trained SVM model to identify the testing set labels (the last panel in Fig.

131 1A).

132 Our DBN-SR based framework can also identify the multi-level temporal features, spatial

133 features, and features for multi-task classification (Fig. 1B). Specifically, the DBN model took

134 fMRI time series from training data as input and produced a weight matrix W for each layer

135 respectively, which represent the multi-layer temporal features of group-wise tfMRI signals

136 (the first two panels in Fig. 1B). These multi-layer temporal features W were served as the

137 temporal dictionary $\boldsymbol{D}^1$ and used as input to the LASSO algorithm to regress corresponding

138 loading coefficient $\boldsymbol{\alpha}^1$, which represents individual-level spatial patterns (the third panel in Fig.

139 1B). Next, the loading coefficient $\boldsymbol{\alpha}^1$ was used as the input of SR stage to derive the common

140 dictionary $\boldsymbol{D}^2$ and the loading coefficient $\boldsymbol{\alpha}^2$, which represent group-wise spatial patterns and

141 features for multi-task classification for each layer, respectively (the last three panels in Fig.

142 1B).

7

Figure 1. The overview of hybrid deep belief network and sparse representation framework (DBN-SR). (A) The pipeline of multi-task fMRI data classification analysis via the proposed model. The seven capital letters refer to seven different tasks respectively (E: emotion, G: gambling, R: relational, M: motor, L: language, S: social, and W: work memory). (B) The detailed illustration of using DBN and SR model to extract multi-level temporal features, spatial features, and features for classification from multi-task fMRI signals. In the second block, the blue line represents temporal features derived from the weights of DBN, while the red line represents task design paradigms.

8

**Data acquisition and preprocessing**

We employed the seven task fMRI data from Q1 release of Human Connectome Project (HCP) in this study (Barch et al., 2013). The details of tfMRI data acquisition and preprocessing pipeline could be referred to our previous study (Song et al., 2022).

Specifically, the seven tasks are emotion, gambling, relational, motor, language, social, and working memory (WM). The number of time points for each task is shown in Table 1. As the tfMRI data consist of different time points, we truncated all tfMRI signals to the same time length (176 frames). In this work, 60 subjects were used from the released dataset

**Table1. Details of the condition and frames for seven tasks**

| TASK | EMOTION | GAMBLING | RELATIONAL | MOTOR | LANGUAGE | SOCIAL | WM |
|------|---------|----------|------------|-------|----------|--------|-----|
| Condition | 2 | 2 | 2 | 6 | 2 | 2 | 8 |
| Frames | 176 | 253 | 232 | 284 | 316 | 274 | 405 |

The truncation preprocessing, unavoidably, influences the integrity of task design. For instance, four conditions are excluded from the WM task due to data truncation. Nonetheless, in terms of other tasks, the truncated tfMRI data include not less than one block for all events (sFig. 1).

**Data preparation**

First, we extracted the whole-brain fMRI signal for each subject using the standard MNI152 template as the mask, resulting in each 2-dimensional matrix. Then the signal matrices of the seven tasks for each subject were spatially concatenated into a large matrix $S_i^1$ ($S_i^1 = [S_{i,E}^1, S_{i,G}^1,$ $S_{i,R}^1, S_{i,M}^1, S_{i,L}^1, S_{i,S}^1, S_{i,W}^1] \in R^{t \times (n \times 7)}$, where $S_{i,E}^1 \in R^{t \times n}$ had $t$ time points and $n$ voxels. The seven capital letter subscripts refer to seven different tasks respectively (E: emotion, G:

9

170  gambling, R: relational, M: motor, L: language, S: social, and W: work memory). TfMRI time

171  series for each voxel were normalized to derive zero mean and unit norm. In this work, five-

172  fold cross-validation scheme was chosen. Thus, 60 subjects were randomly divided into five

173  equal folds. In each iteration, one fold (12 subjects) was taken for testing and the rest four (48

174  subjects) for training. It is noteworthy that the training and testing sets for each iteration were

175  completely independent. Then, the multi-task fMRI signal matrices of all the subjects in the

176  training set were spatially concatenated to compose a multi-subject fMRI matrix $\boldsymbol{S}^1 = [\boldsymbol{S}_1^1,$

177  $\boldsymbol{S}_2^1, ..., \boldsymbol{S}_p^1] \in R^{\text{t} \times (\text{n} \times 7 \times p)}$, where $p$ is the number of training subjects ($p = 48$) (Fig. 1A).

178     As whole-brain fMRI data generally contain enormous voxels, the group-wise tfMRI

179  signals consisting of multiple tasks and subjects exhibit relatively high dimensionality,

180  inevitably resulting in an overloaded computational burden and memory consumption. To

181  tackle these problems, we randomly selected only 10% of voxels' whole-brain signals for each

182  subject in training stage (Huan Liu 2017; Song et al., 2022). To ensure the uniform distribution

183  of sampled voxels across different brain regions, we employed the Fisher-Yates shuffle

184  algorithm implemented by the "randperm" function in MATLAB, known for generating

185  random permutations with a uniform distribution (Fisher & Yates, 1938). The distribution of

186  the randomly selected 10% voxels across all subjects can be found in the Supplementary

187  Materials (sFig. 6-7).


188  **Deep belief network model-based analysis**


189  In this work, we chose DBN to extract group-wise temporal features based on previous research

190  demonstrating its ability to identify meaningful FBNs (Qiang et al., 2020; Ren et al., 2021). In

10

191   general, DBN can be regarded as stacked blocks of Restricted Boltzmann Machines (RBM) (G.

192   E. Hinton, Osindero, & Teh, 2006), an energy-based probability generation model that

193   simulates the potential distribution of input data via interactions between visible and hidden

194   variables. While units between visible layer $v$ and hidden layer $h$ are connected by weights,

195   there is no connection within the layer. As a multiple stacked RBM model, the DBN model is

196   designed to learn and train weights for each layer. As described in Asja Fischer (2012) and X.

197   Hu et al. (2018), the energy function of the DBN model adopted to update the weights layer by

198   layer is defined as follows:

$$E(v,h) = \sum b_i v_i - \sum b_j h_j - \sum v_j h_j w_j \tag{1}$$

200   Where $v_i$ and $h_j$ refer to the activation state of two layers; $b_i$ and $b_j$ represent their bias; $w_j$

201   indicate the weight between layer $i$ and layer $j$.

202      As introduced in the previous section, the tfMRI signals of randomly selected 10% voxels

203   in each individual's whole brain of multi-task in training set were spatially concatenated to

204   generate a multi-subject fMRI matrix for model training, and thus the group-wise tfMRI time

205   series (176 time points) were taken as training samples for the DBN model. In our work, the

206   neural architecture of DBN model was set as 4 layers and 128 neurons experimentally and

207   empirically (see Parameter Selection part). Specifically, the number of visible variables $t$ is the

208   same as the number of time points of fMRI signal (i.e., 176 in our study), and the number of

209   hidden variables $k1$ in each hidden layer represents the number of latent components expressed

210   in fMRI data ($k1$=128). The DBN model was adopted to model group-wise tfMRI matrix $S^1$

211   to obtain a weight matrix $w_j$ from each layer. The weight matrix of visible layer is represented

212   by $w_1 \epsilon R^{t \times k1}$, and the weight matrix of each hidden layer refers to $w_j \epsilon R^{k1 \times k1}$ ($j$ =2,3,4). The

11

213　multi-layer temporal features $W_j$ in each layer of DBN model can be derived by successive

214　multiplication of the weight matrices on the adjacent layers ( $W_j \epsilon R^{t \times k1}$ ), that is,

215　$W_4 = w_4 * w_3 * w_2 * w_1$, $W_3 = w_3 * w_2 * w_1$, $W_2 = w_2 * w_1$ , $W_1 = w_1$. Since each sample

216　input to the DBN model consists of all time points for each voxel, the weights $w_j$ (j =1,2,3,4)

217　across 4 layers represent the temporal features of the input fMRI data at different levels of

218　abstraction. Thus, the successive multiplication of weight matrix $W_j$ (j =1,2,3,4) obtained from

219　each layer of the DBN model represents multi-level temporal features embedded in fMRI

220　signals.

221　　　Drawing inspiration from the successful application of LASSO regression for deriving

222　spatial features in previous studies (Haufe et al., 2014; Lee, Jeong, & Ye, 2013), we performed

223　the LASSO regression to derive individual spatial features. Specifically, the multi-layer

224　temporal features $W_j$ derived by the DBN model were normalized and then served as the

225　temporal dictionary $D^1 \epsilon R^{t \times k1}$ (Calhoun et al., 2001; Tibshirani, 2011). Here, as the successive

226　multiplication of weight matrices leads to the larger scale of deeper dictionaries, a

227　normalization procedure ensures reasonable performance of LASSO regression at the same

228　scale. Subsequently, we employed the original individual signal matrix $S_i(i \in 1, 2, …, p)$,

229　along with the temporal dictionary $D^1$ as input to the LASSO algorithm, which produce the

230　corresponding individual loading coefficient $\alpha_i^1$ ($\alpha_i^1 \in R^{k1 \times n}$, n=228453). Since $D^1$

231　incorporates the group-wise temporal features, the resulting individual loading coefficients $\alpha_i^1$

232　obtained through regression can be considered as spatial sparse representations of each

233　individual's fMRI signals $S_i$ on the common temporal dictionary $D^1$. Consequently, the

234　individual loading coefficients $\alpha_i^1$ represent the individual spatial features. Here, all the loading

12

235    coefficient matrix derived from LASSO regression refers to $\boldsymbol{\alpha}^1$ ($\boldsymbol{\alpha}^1=[\boldsymbol{\alpha}_1^1, \boldsymbol{\alpha}_2^1, ..., \boldsymbol{\alpha}_i^1, ..., \boldsymbol{\alpha}_p^1]$

236    $\in R^{\,k1\times(n\times7\times p)}$, $\boldsymbol{\alpha}_i^1 = [\boldsymbol{\alpha}_{i,E}^1, \boldsymbol{\alpha}_{i,G}^1, \boldsymbol{\alpha}_{i,R}^1, \boldsymbol{\alpha}_{i,M}^1, \boldsymbol{\alpha}_{i,L}^1, \boldsymbol{\alpha}_{i,S}^1, \boldsymbol{\alpha}_{i,W}^1] \in R^{k1\times(n\times7)}$.

237    Similarly, in order to derive the loading coefficient matrix $\boldsymbol{\alpha}_{test}^1$ for testing set of each

238    layer, the group-wise time-series dictionary matrix $\boldsymbol{D}^1$ derived from the training stage was

239    applied to model $\boldsymbol{S}_{test}^1$ to obtain $\boldsymbol{\alpha}_{test}^1$ by resolving a typical l-1 regularized LASSO problem.

240    In this work, the regularization parameter $\lambda 1$ of LASSO regression was set as 0.1

241    experimentally and empirically.

**Sparse Representation model**

243    Although we successfully obtained individual loading coefficient matrices $\boldsymbol{\alpha}^1$ and $\boldsymbol{\alpha}_{test}^1$

244    through LASSO regression for the training and testing sets, respectively, these features were

245    unsuitable for classification due to their high dimensionality ($\boldsymbol{\alpha}^1 \in R^{k1\times n}$, $k1=128$, $n=228453$).

246    Therefore, our next goal was to extract the multi-level group-wise spatial patterns based on the

247    individual spatial patterns, and finally excavate multi-level features for multi-task classification

248    that could distinguish multi-task fMRI signals and reveal the distinctive organization patterns

249    of different task stimulations. Here, we adopted a sparse representation based model, which

250    has already been proven as an effective algorithm in previous research to identify the intrinsic

251    spatial functional patterns and features for multi-task classification from fMRI data (Song et

252    al., 2022; S. Zhang et al., 2016). Specifically, we first aggregated all the loading coefficient

253    matrices $\boldsymbol{\alpha}_i^1$ of all the subjects into one matrix $\boldsymbol{S}^2$ for each layer of the DBN model ($\boldsymbol{S}^2 = [\boldsymbol{S}_1^2,$

254    $\boldsymbol{S}_2^2,...,\boldsymbol{S}_i^2,...,\boldsymbol{S}_p^2] \in R^{k1\times(n\times7\times p)}$, where $\boldsymbol{S}_i^2 = [(\boldsymbol{\alpha}_{i,E}^1)^{\mathrm{T}}, (\boldsymbol{\alpha}_{i,G}^1)^{\mathrm{T}}, (\boldsymbol{\alpha}_{i,R}^1)^{\mathrm{T}}, (\boldsymbol{\alpha}_{i,M}^1)^{\mathrm{T}}, (\boldsymbol{\alpha}_{i,L}^1)^{\mathrm{T}}, (\boldsymbol{\alpha}_{i,S}^1)^{\mathrm{T}},$

255    $(\boldsymbol{\alpha}_{i,W}^1)^{\mathrm{T}}] \in R^{n\times(7\times k1)}$. Then, $\boldsymbol{S}^2$ would be served as the input for dictionary learning and sparse

13

256      representation to derive a group-wise spatial dictionary $D^2 \in R^{n \times k2}$ and the corresponding

257      loading coefficients $\alpha^2$ for each layer, respectively. Note that $k2$ represents the number of

258      dictionary atoms, which was set as the same value as $k1$ ($k2=128$). Here, $\alpha^2=[\alpha_1^2, \alpha_2^2, \ldots,$

259      $\alpha_i^2, \ldots, \alpha_p^2] \in R^{k2 \times (k1 \times 7 \times p)}$, where $\alpha_i^2=[\alpha_{i,E}^2, \alpha_{i,G}^2, \alpha_{i,R}^2, \alpha_{i,M}^2, \alpha_{i,L}^2, \alpha_{i,S}^2, \alpha_{i,W}^2] \in R^{k2 \times k1 \times 7}$.

260      The loss function of sparse representation model yields a sparse resolution constraint on the

261      loading coefficient $\alpha^2$ with an l1 regularization (Eq. (2)), where $\lambda2$ is a regularization

262      parameter that can balance the regression residual and sparsity level. $\lambda2$ was set as 0.05.

$$Min\frac{1}{2}\|S^2 - D^2\alpha^2\|_F^2 + \lambda2\|\alpha^2\|_{1,1} \qquad (2)$$

264      To prevent $D^2$ from arbitrarily large values that cause the trivial solution of the

265      optimization, the columns $d_1, d_2, \ldots, d_k$ are restricted by Equation (3).

$$C \triangleq \{D^2 \in R^{t \times k2}, s.t. \forall j = 1, \cdots, k2, \quad d_j^T d_j \leq 1\} \qquad (3)$$

267      As the dictionary $D^2$ was obtained by a sparse representation of $\alpha^1$, which comprise all

268      individual spatial features, the learned dictionary $D^2$ consequently represents the group-wise

269      spatial features. Correspondingly, $\alpha_i^2$ was a sparse representation on the common spatial

270      dictionary $D^2$. Given the ability of a sparse representation model to effectively reduce the

271      dimensionality of raw fMRI data while retaining its essential information, the resulting intrinsic

272      features ($\alpha_i^2$) derived from the extraction of common temporal and spatial dictionaries can

273      effectively capture the variations in spatio-temporal patterns of functional brain activity across

274      different tasks. As a result, these intrinsic features were suitable for multi-task classification.

275      To derive the $\alpha_{test}^2$ of testing set for post-hoc classification analysis, we also leveraged

276      the LASSO regression algorithm for each layer. Specifically, the loading coefficient matrix

277      $\alpha_{test}^1$ was regarded as the input matrix $S_{test}^2$, and the dictionary matrix $D^2$ derived from the

278    training stage was employed to model $S_{test}^2$ to learn the loading coefficient $\alpha_{test}^2$. All the

279    parameters in testing stage were set the same as in the training stage.

**Parameter Selection**

281    The determination of hyperparameters, such as the number of cross-validation folds, the

282    number of layers and neurons of the DBN model, and the regularization parameters of the

283    sparse representation model, was accomplished through a combination of referring to previous

284    studies and learning from the training set, the testing set was not involved in any parameter

285    selection process.

286        The choice of cross-validation folds is crucial as it offers a trade-off between precision

287    and computational cost for performance estimation (Hansen et al., 2013). Commonly used

288    cross-validation folds in current machine learning experiments often include 2-fold, 5-fold, 10-

289    fold, or the leave-one-out method. In theory, while some studies suggest the 10-fold or leave-

290    one-out method may provide a higher estimated accuracy (Kohavi, 1995), some reveals that 5-

291    fold or 10-fold is the optimal choice for balancing computational cost and accuracy (Hansen et

292    al., 2013). However, due to the need for our framework to combine all individuals within the

293    training set to extract group-wise temporal features during training phase, the computational

294    resource demands of the 10-fold or leave-one-out method are greater. Therefore, we opted for

295    the 5-fold approach. To further validate our selection, we conducted a comparative analysis

296    between the 2-fold and 5-fold to assess the decoding accuracy. The findings revealed that the

297    average decoding rate was slightly lower for the 2-fold compared to the 5-fold, providing

298    additional confirmation of our initial selection. (sTab. 1).

299    Our selection of a 4-layer, 128-neuron DBN structure was set based on our previous study

300    utilizing the neural architecture search technique (NAS) for recognizing spatio-temporal

301    features from fMRI data (Xu, Ren, Tao, Song, & He, 2022),which effectively determined the

302    optimal structure for DBN model with 3 layers and 120-150 neurons. Therefore, in our study,

303    we defined the number of neurons as 128 and experimented with both 3-layer and 4-layer

304    configurations to extract meaningful task-related temporal features. Specifically, we compared

305    the group-wise temporal features derived from DBN model with 3-layer and 4-layer structures,

306    in terms of their Pearson correlation coefficient (PCC) with task paradigm curve, based on

307    training set (fold 5). The results revealed that the 4-layer DBN outperformed in capturing

308    temporal features, as indicated by the higher PCC values observed in 4-layer structure (Tab. 2).

309    In terms of selecting the number of neurons, we took into consideration computational

310    efficiency. We determined that selecting 128 neurons, a power of two within the desired range

311    of 120-150, would optimize computational speed. Hence, we concluded that the optimal

312    configuration for the DBN model with 128 neurons and 4 layers.

313    The regularization parameter ($\lambda$) plays a crucial role in sparse representation and LASSO

314    regression. Although no golden standard exists for determining the value of $\lambda$, previous studies

315    on FBN recognition have experimentally set $\lambda$ within the range of 0.05 to 0.5 (Fangfei Ge,

316    2018; Lv, Jiang, Li, Zhu, Chen, et al., 2015; Shu Zhang 2017). In our previous work on task

317    fMRI data classification using a two-stage sparse representation approach, we conducted

318    parameter selection experiments within the range of $\lambda$ from 0.05 to 0.5 and found that the

319    highest accuracy was achieved when $\lambda 1 = 0.1$ and $\lambda 2 = 0.05$ or 0.1 (Song et al., 2022). Here, $\lambda 1$

320    and $\lambda 2$ represent the regularization parameters for the LASSO regression and sparse

16

321 representation, respectively. Therefore, in this study, we determined the λ1 as 0.1, and

322 systematically changed the setting of the regularization parameter in the sparse representation

323 λ2 (λ2=0.05, 0.1) while evaluating their impact on the obtained group-wise spatial features

324 derived from training set (fold 5). The results showed that when λ2 was set to 0.05, a greater

325 number of FBNs could be identified in the group-wise spatial features $\boldsymbol{D}^2$ by comparison with

326 the general linear model (GLM) -derived activation patterns (Tab. 3). Consequently, we set

327 λ1=0.1 and λ2=0.05 as regularization parameters for LASSO regression and sparse

328 representation stage, respectively. To further validate this, we assessed the classification

329 accuracy on testing dataset using these two different λ2 values (0.05, 0.1) while keeping λ1=0.1

330 for all 5 folds. The results demonstrated that λ2=0.05 achieved higher accuracy, reconfirming

331 our choice (sTab. 2).

332 **Table 2. Comparison of Pearson correlation coefficient (PCC) for 3-layer structure and**

333 **4-layer structure.**

| Structure | Layer1 | Layer2 | Layer3 | Layer4 | Mean±SD |
|-----------|--------|--------|--------|--------|---------|
| 3-layer | 0.48±0.12 | 0.52±0.06 | 0.50±0.06 | | 0.50±0.08 |
| 4-layer | 0.55±0.00 | 0.63±0.01 | 0.66±0.03 | 0.71±0.02 | 0.64±0.02 |

334 **Table 3. Comparison of the number of identified FBNs cross each layer for different λ2**

335 **values.**

| λ2 | Layer1 | Layer2 | Layer3 | Layer4 |
|------|--------|--------|--------|--------|
| 0.05 | 15 | 17 | 22 | 45 |
| 0.1 | 12 | 13 | 18 | 27 |

336 **Identification of multi-level temporal patterns**

337 As mentioned in the "Deep belief network model based analysis" section, $W_j$ of the $j$-th hidden

338 layer ($j = 1,2,3,4$) represents the temporal features of group-wise tfMRI for respective layer

17

339   (Fig. 1B). Here we used PCC as a metric to identify the task-related temporal features (Benesty,

340   Chen, Huang, & Cohen, 2009; Lv, Jiang, Li, Zhu, Chen, et al., 2015). Specifically, we first

341   calculated the task paradigm curves convolved with hemodynamic response function (HRF).

342   Next, we computed the PCC values between the convolved task paradigm curves and the atoms

343   in the group-wise temporal features $\boldsymbol{D}^1$ derived from the DBN model, following standard

344   procedures employed in previous studies (Kay, Rokem, Winawer, Dougherty, & Wandell, 2013;

345   O'Reilly, Woolrich, Behrens, Smith, & Johansen-Berg, 2012). The PCC of the identified

346   temporal features and the task-based stimulus can be defined as Equation (4).

347 $$P_{corr, \, c} = corr \, (\boldsymbol{D}_c^1, \, TASK) \qquad (4)$$

348   Here, $\boldsymbol{D}_c^1$ refers to the c-th component in temporal features $\boldsymbol{D}^1$ derived from DBN stage (c = 1,

349   $\cdots, k$ 1). TASK represents the task paradigm curves convolved with HRF. Essentially, $P_{corr, \, c}$,

350   measures the temporal similarity between the temporal patterns of $\boldsymbol{D}_c^1$ and the task stimulus.

351   The atoms with the highest PCC value in group-wise temporal features $\boldsymbol{D}^1$ were chosen to

352   represent the multi-layer temporal features.

**Identification of multi-level spatial patterns**

353

354   The multi-level spatial patterns can also be identified in the second stage of sparse

355   representation model. Specifically, the $\boldsymbol{S}_{i,t}^1$ can be factorized into $\boldsymbol{D}^1$ and the loading

356   coefficient $\boldsymbol{\alpha}_{i,t}^1$, which represent the group-wise temporal features and the individual spatial

357   features, respectively. Here, $i$ refers to $i$-th subjects (i ∈ 1, 2, …, p, and p=48 in this work), $t$

358   means $t$ kind of task, $t \in \boldsymbol{\Phi} = \{E, G, R, M, L, S, W\}$. To further derive the group-wise spatial

359   features, the transposition of $\boldsymbol{\alpha}^1$ could be then decomposed into $\boldsymbol{D}^2$ and $\boldsymbol{\alpha}^2$ as shown in

18

360  Equation (5). Since the transpose of $\boldsymbol{\alpha}_{i,t}^1$ can be expressed as dictionary $\boldsymbol{D}^2$ multiplied by

361  loading coefficient $\boldsymbol{\alpha}_{i,t}^2$ (Equation (5)), the relationship between $\boldsymbol{S}_{i,t}^1$ and $\boldsymbol{D}^1$, $\boldsymbol{D}^2$, $\boldsymbol{\alpha}^2$ can be

362  deduced as Equation (6) shown, which also consistent with previous studies (Huan Liu 2017;

363  Song et al., 2022).

$$\boldsymbol{S}_{i,t}^2 = (\boldsymbol{\alpha}_{i,y}^1)^T = \boldsymbol{D}^2 \times \boldsymbol{\alpha}_{i,t}^2 \tag{5}$$

$$\boldsymbol{S}_{i,t}^1 = \boldsymbol{D}^1 \times \boldsymbol{\alpha}_{i,t}^1 = \boldsymbol{D}^1 \times (\boldsymbol{D}^2 \times \boldsymbol{\alpha}_{i,t}^2)^T \tag{6}$$

366  Since all subjects share the same group-wise temporal dictionary $\boldsymbol{D}^1$, the common

367  dictionary $\boldsymbol{D}^2$ contained group-wise spatial patterns, of which atoms could be used to define

368  the FBNs. Thus, the corresponding multi-layer spatial features were derived from the common

369  dictionary $\boldsymbol{D}^2$ for each layer of the proposed framework (the fourth and fifth panels in Fig. 1B).

370  We then identified the spatial correlation coefficient (SCC) to quantify the similarity

371  between spatial patterns obtained from the proposed framework and the GLM -derived

372  activation patterns. Specifically, the GLM-based analysis was performed individually, followed

373  by group-wisely analysis using FSL FEAT (http://www.fmrib.ox.ac.uk/fsl/feat5/index.html).

374  The group-level GLM-based results were employed for comparison. More details of GLM

375  analysis are available in previous literature (Lv, Jiang, Li, Zhu, Zhang, et al., 2015). The SCC

376  is defined in Equation (7) (Ben J. Harrison, 2008; Zuo et al., 2010):

$$\mathbf{R}\ (\boldsymbol{X}, \boldsymbol{T}) = \frac{\Sigma_{p=1}^n (X_p - \bar{X})(T_p - \bar{T})}{\sqrt{\Sigma_{p=1}^n (X_p - \bar{X})^2 \cdot \Sigma_{p=1}^n (T_p - \bar{T})^2}} \tag{7}$$

378  where $\boldsymbol{X}$ is the spatial functional network derived by the proposed framework, $\boldsymbol{T}$ represents

379  the GLM-derived activation template, and $n$ refers to the number of voxels of whole brain.

19

380 **SVM-based classification method**

381 To further classify multi-task fMRI signals, we performed five-fold cross-validation to evaluate

382 the classification performance of the proposed framework. As the linear SVM has optimization

383 and generalization capability in limited sample sizes, as well as its proven effectiveness in

384 multi-class classification (Chang & Lin, 2011b; Jang et al., 2017), we conducted multi-task

385 classification analysis based on linear SVM classifier, which was established by the LIBSVM

386 toolbox (Chang & Lin, 2011a). For each layer, as the loading coefficient $\boldsymbol{\alpha}^2$ contains both

387 temporal and spatial features embedded in fMRI signals, we first trained the SVM classifier

388 using $\boldsymbol{\alpha}^2$ derived from training set, and then evaluated the classification performance by

389 feeding the $\boldsymbol{\alpha}^2_{test}$ of testing set into the trained SVM model. Based on the true label of seven

390 tasks for each loading coefficient $\boldsymbol{\alpha}^2_{test}$, the classification accuracy of each layer in each fold

391 was defined as the percentage of correctly predicted samples. The final classification accuracy

392 for each layer is the average of five folds for seven tasks. We then calculated the specificity of

393 each fold for each layer, and the final specificity for each layer is the average of the five folds.

394 **ROA-based analysis**

395 The further goal aimed at uncovering discriminative functional components for multi-task

396 classification. Inspired by the successful use of the Ratio of activation (ROA) in identifying

397 discriminative components for decoding resting state fMRI (rsfMRI) and tfMRI (S. Zhang et

398 al., 2016), we raised a novel ROA metric to identify the key components for seven-task

399 classification. The ROA of the $j$-th row in loading coefficients $\boldsymbol{\alpha}^2$ could be defined as follows:

$$N_t = |\boldsymbol{\alpha}^2(j,k)|_0, kth \ column \ belongs \ to \ task(t)$$

20

$$401 \qquad \mathrm{ROA}_j = \sqrt{\frac{1}{T}\sum_{t=1}^{T}(N_t - \overline{N_t})^2} \qquad (8)$$

402       In Equation (8), $\boldsymbol{\alpha}^2$ represent all the individual spatio-temporal features, $\boldsymbol{\alpha}^2 = [\boldsymbol{\alpha}_1^2, \boldsymbol{\alpha}_2^2, \ldots,$

403       $\boldsymbol{\alpha}_i^2, \ldots, \boldsymbol{\alpha}_p^2] \in R^{k2 \times (k1 \times 7 \times p)}$ ($k1 = k2 = 128$, p=48). $i$ refers to $i$-th subject ($i \in 1, 2, \ldots,$ p). $t$

404       represents task index (t$\in$1, 2, ..., 7), and $T$ represents the number of task paradigms (i.e., 7 in

405       our work). Task ($t$) represents each of the seven different tasks. $N_t$ represents the activation

406       level for each task, and $\overline{N_t}$ represents the average of $N_t$ ($t = 1, \cdots, 7$). Here, the activation level

407       $N_t$ was defined by counting the number of non-zero entries marked as each task in the

408       corresponding each row vector of $\boldsymbol{\alpha}^2$ (t$\in$1, 2, ..., 7). As $\boldsymbol{\alpha}^2$ is a sparse matrix, the task with a

409       higher count of nonzero elements in the row vectors of $\boldsymbol{\alpha}^2$ is deemed to be more "active".

410       Therefore, $N_t$ represents each task's activation level in the row vectors of $\boldsymbol{\alpha}^2$. The ROA was

411       calculated by counting the standard deviation of $N_t$ across the seven tasks. A larger ROA value

412       (i.e., larger standard deviation) indicates greater differences in activity levels across the seven

413       tfMRI signals, which were more discriminative for multi-task classification.

414       To validate that the components of higher ROA values capture greater capacity in

415       classifying the multi-task fMRI signals, an experiment was designed as illustrated below. After

416       sorting the ROA values for all components (i.e., rows in loading coefficients $\boldsymbol{\alpha}^2$) from highest

417       to lowest, we iteratively adopted more rows sorted by their ROA values in $\boldsymbol{\alpha}^2$ as feature inputs

418       for training the SVM classifier, that is, the components with higher ROA values were used

419       preferentially for training. Afterwards, the corresponding components of $\boldsymbol{\alpha}_{test}^2$ from testing set

420       were entered into the trained SVM model to evaluate the classification accuracy. Specifically,

421       to define the key components with greater capacity for multi-task classification in each layer,

422       we have repeated this ROA analysis using $\boldsymbol{\alpha}^2$ derived from each layer of proposed model. Here

21

423 we applied the same classification scheme described in the previous section "SVM-based

424 classification method".

425      After establishing the ROA metric for the classification features $\boldsymbol{\alpha}^2$, our subsequent

426 objective is to elucidate the neural implications of these classification features. Given that each

427 row of $\boldsymbol{\alpha}^2$ corresponds to each column of $\boldsymbol{D}^2$ (i.e., each atom in $\boldsymbol{D}^2$), and these atoms can be

428 mapped back to brain space, we thus established a relationship between the brain activations

429 derived from the atoms in $\boldsymbol{D}^2$ and the ROA values of the row vectors of $\boldsymbol{\alpha}^2$. This connection

430 allows us to interpret neural implications of classification features.

431 **Result**

432 **Classification performance of multi-task fMRI signals**

433 By applying the proposed DBN-SR framework to multi-task fMRI data using five-fold cross-

434 validation strategy, our results reveal that the fMRI data of seven tasks can be accurately

435 classified. In detail, the classification accuracy for five-fold ranges from 92.86% to 100%, with

436 an average accuracy of 97.86%±3.42% (Mean ± SD) in the layer 4 (Fig. 2A), which

437 demonstrated the proposed framework can effectively uncover the inherent differences in

438 composition patterns of multi-task fMRI signals.

439      We also explored the classification performance based on features derived from each layer

440 of the proposed framework (Fig. 2). The trend of the classification accuracy curves for five

441 folds is relatively steady, with an average accuracy of 98.15%±0.90% (Mean±SD) (Fig. 2A).

442 Moreover, the average accuracies across five-fold from layer1 to layer4 are 99.29%, 98.33%,

443 97.14%, and 97.86%, respectively. We depicted confusion matrices for each layer to represent

444 the average classification accuracy of the seven tasks, as shown in Figure 2b. The results

445 indicate that all the average classification accuracies for seven tasks across five-fold are greater

446 than 95% in each layer, except for three major confusions, that is, gambling task in layer 3 and

447 layer 4, relational task in layer 2 and layer 3, and language task in layer 3 (Fig. 2B). In addition,

448 the specificity of classification results of the first two layers is slightly higher than that of the

449 deeper two layers (Fig. 2C). Overall, the classification performance of the shallower layers is

450 relatively better than that of the deeper layers.

451

452 Figure 2. Classification performance. (A) The classification accuracy of five-fold in each layer.

453 (B) The average confusion matrices of five-fold cross-validation on the seven tasks. (C) The

454     average specificity of five-fold cross-validation classification on the seven tasks.

455     **Identified multi-level temporal and spatial patterns of multi-task fMRI signals**

456     Multi-level temporal patterns

457     Our DBN-SR based framework can effectively identify the temporal patterns of multi-task

458     fMRI signals at multi-scale (Fig. 3). In each layer, we quantitatively compared the PCC of the

459     identified temporal features and each task-based stimulus. Those atoms with the highest PCC

460     value in temporal dictionary $D^1$ were chosen to represent the task-related temporal patterns.

461     We randomly select one training fold as an example to show the representative temporal

462     patterns for each layer (fold 5) (Fig. 3). The average PCC values of seven tasks for all 5-fold

463     can be found in Supplemental Table 6.

464        The overall multi-level temporal patterns are relatively consistent with the task design

465     paradigms. Specifically, the average PCC of seven tasks from layer1 to layer4 is 0.55±0.12,

466     0.61±0.03, 0.65±0.07, and 0.71±0.08 (Mean ± SD), respectively, where the highest correlation

467     is observed in layer4 (Fig. 3). Intriguingly, there exist gradient in the resolution of temporal

468     patterns derived from different layers. In the shallow layer, all the identified temporal patterns

469     are mixed with many random noises, resulting in a relatively poor correlation with task

470     paradigms. In comparison, in the deeper layer, the temporal patterns are smoother and more

471     consistent with the original task design curves, indicating that DBN-SR model can filter noises

472     in each layer while keeping useful information of brain activities, which agrees with the former

473     research (H. Huang et al., 2018; Wei Zhang, 2020).

Figure 3. Comparison of group-wise temporal patterns for seven tasks across different layers, including the identified temporal features (blue lines) and the task paradigms (red lines). The quantitative similarities (PCC) of identified temporal features with task paradigms are also provided. The y-axis represents the stimulus response amplitude, while the x-axis represents time point. The background colors represent different layers of our DBN-SR model. The lighter colors represent shallower layers, while the darker colors represent deeper layers.

Multi-level spatial patterns

Our framework can also effectively identify the spatial patterns from different layers. The most predominant spatial patterns identified by the proposed framework are the task-evoked FBNs, including emotion, gambling, relational, motor, social, language, and working memory. In each layer, we quantitatively compared the SCC of the identified spatial patterns and the GLM-derived activation patterns. Those atoms with the highest SCC value in spatial dictionaries $D^2$ were chosen to represent the spatial pattern. We randomly selected one training fold to illustrate

488    the representative FBNs for each layer (Fig. 4).

489        Overall, the spatial patterns are generally consistent with the GLM-derived activation

490    patterns, with increasingly precise resolution from shallow to deep layers. Quantitatively, the

491    average SCC of seven tasks from layer1 to layer4 is 0.36±0.20, 0.26±0.11, 0.40±0.12, and

492    0.48±0.12 (Mean ± SD), respectively, where the highest SCC is observed in layer 4 (Fig. 4).

493    Intriguingly, there exist distinct differences among spatial patterns derived from different layers.

494    The spatial patterns across layers show a trend of increasing consistency with the GLM-derived

495    activation patterns, and are more compact in deeper layers for most tasks. Meanwhile, more

496    FBNs can be found in the deeper layers compared with shallow layer. For example, some FBNs

497    cannot be identified in the first three layers, such as FBNs related to gambling and relational

498    tasks (Fig. 4).



Figure 4.

501    Comparison of group-wise spatial patterns for seven tasks across different layers. The spatial

502    correlation coefficient (SCC) between each identified spatial pattern and GLM-derived

503    activation pattern is labeled on top of each brain map.

504        Apart from FBNs, the proposed framework can also effectively detect various artifact-

505    related components. Specifically, the atoms in spatial dictionary $D^2$ can represent the group-

26

506     wise spatial features and can be mapped back to the 3D brain volume. Subsequently, we

507     manually inspected whether spatial map matched the known types of artifacts based on

508     previous study (Salimi-Khorshidi et al., 2014). Through this process, we found several artifact-

509     related components, including movement-related, cardiac-related, sagittal sinus, susceptibility-

510     motion, white-matter, and MRI acquisition/reconstruction related (Fig. 5).



511

512     Figure 5. Identified artifact components, including movement-related, cardiac-related, sagittal

513     sinus, susceptibility-motion, white-matter, and MRI acquisition/reconstruction related.

514     Overall, our effective DBN-SR model is capable of characterizing the multi-level

515     spatiotemporal features of brain function. The quantitative analysis further demonstrates that,

516     in deeper layer, the representative temporal features correspond well with task design curves,

517     and the spatial features are relatively more consistent with the GLM-derived activation. In

518     addition to task-evoked functional components, our framework could also effectively identify

519     artifact components from group-wise multi-task fMRI data, laying the groundwork for further

520     research into the functional role of these components in multi-task classification.

521     **Identification of discriminative features by ROA analysis**

522     As depicted in the "ROA-based analysis" section, we first computed the ROA index by sorting

523     the ROA values of all the components in loading coefficients $\alpha^2$ of the training set, then, in

524     order to evaluate the classification performance, the corresponding components in the loading

27

525     coefficient $\boldsymbol{\alpha}^2_{test}$ of testing set were fed sequentially into the trained SVM classifier according

526     to the ROA index. Here, the classification results of each layer on one randomly selected testing

527     fold dataset (fold 5) using different number of components, sorted by their ROA values, are

528     illustrated in Fig. 6A. While the number of components increases from 1 to 20, the accuracy

529     curves of four layers grow monotonically, and the average accuracy of all curves rises to

530     91.96%. When more than twenty components are included for classification, the accuracy

531     curves of four layers exhibit a plateau with accuracies reaching close to 100%, indicating that

532     the additional components with lower ROA values contribute less to the successful

533     classification of multi-task signals. Thus, the top twenty components with higher ROA values

534     can be regarded as key components for the classification task to some extent. Generally, our

535     method can effectively disclose the key components with great classification capacity. In

536     addition, the findings are consistent across different testing folds, hence the additional results

537     of the other four folds are included in the Supplementary Materials (sFig2-5).

538         To further investigate the neural implications of key components with greater

539     classification capacity, we inspected the spatial patterns of the top twenty key components

540     identified by ROA analysis in each layer. By further analyzing the composition of the twenty

541     key components in each layer, we found that these key atoms are either FBNs or artifact-related

542     components, which were identified by visually examining their spatial patterns with established

543     templates and further calculating their SCC with GLM-derived activation maps.

544         Intriguingly, our results show that the top twenty key components in the four layers are

545     largely composed of artifacts, while the proportion of FBNs in key components is small as a

546     whole. On the other hand, the proportion of FBNs is relatively higher in deeper layers compared

28

547 to shallower layers (Fig. 6B). This conclusion aligns with the findings when using the top 40

548 components as key components (sFig. 8).



549

550 Figure 6. ROA classification results in each layer (fold 5). (A) Classification accuracy for

551 SVM-based classification of four layers using the different number of components sorted by

552 their ROA values. (B) The composition of twenty key components sorted by ROA value across

553 each layer.

## Discussion

555 In this study, we proposed a hybrid spatio-temporal deep belief network and sparse

556 representation framework to decode multi-task fMRI signals on a relatively small cohort

29

557 dataset. Our framework could classify fMRI signals of seven tasks with high accuracy and

558 detect multi-level temporal patterns and FBNs, suggesting the effectiveness of the proposed

559 method. In addition, our framework can reveal key components including artifact components

560 and functional brain networks in multi-task classification and uncover their underlying

561 neurological implication.

562     Our proposed framework is composed of several elements, including DBN model,

563 LASSO regression, sparse representation, and SVM classifier, resulting in a relatively complex

564 structure. Nevertheless, our framework achieved a relatively higher classification accuracy in

565 comparison to prior research that also conducted classification of 7 task states on the HCP

566 dataset (X. Huang, Xiao, & Wu, 2021; Wang et al., 2020), while also yielding interpretable

567 classification components. Specifically, Wang et al. (2020) reported two standard machine

568 learning algorithms, namely MVPA-SVM and DNN, and X. Huang et al. (2021) proposed a

569 novel framework (CRNN) incorporating multiple modules such as CNN, recurrent neural

570 network (RNN), and attention mechanism. The average accuracy of our framework (98.15%)

571 is much higher than that of MVPA-SVM (69.2%) and comparable to the accuracies of DNN-

572 based model (93.7%) and CRNN-based model (94.31%) (X. Huang et al., 2021; Wang et al.,

573 2020). Additionally, the neuroscientific implications of their results remain elusive. In

574 conclusion, our proposed model achieved higher decoding accuracy than these models, while

575 also providing a more comprehensive and interpretable methodology for decoding fMRI data.

576     Furthermore, our model unveils multi-level temporal and spatial patterns, demonstrating

577 a resolution gradient spanning from shallow to deep layers. Specifically, in the deeper layers,

578 the identified temporal features are better correlated to the original task paradigm curves.

579     Meanwhile, more diverse FBNs can be detected and the spatial features show more consistency

580     with the GLM-derived activation patterns, in deeper layers.

581     Intriguingly, although more higher-order FBNs can be detected in deeper layers, the

582     classification accuracy using features for multi-task classification derived from deeper layers

583     is lower than that of shallower layers, indicating that these higher-order FBNs are not much

584     helpful for multi-task classification. To validate this observation, we specifically selected only

585     FBNs components from all available components across all five folds for multi-task

586     classification, resulting in an average accuracy of 97.08%±2.14% (Mean±SD), slightly lower

587     than the classification rate obtained using all components (98.15%±0.90%) (sTab. 3). The

588     possible reason is that the FBNs evoked by different cognitive tasks may have co-activated

589     brain regions, thus the FBNs components alone may not fully reveal the potential fundamental

590     differences in functional composition patterns of multi-task fMRI data. On the other hand,

591     ROA-based analyses indicate that artifact components occupy higher proportion of key

592     components for multi-task classification in shallower layers than that in deeper layers, along

593     with higher classification accuracy and specificity in the shallower layers. These findings

594     suggest that the artifact components play an important role in multi-task fMRI signal

595     classification, which is also consistent with previous research, where the artifact components

596     of the EEG signal are significantly more informative than brain activity concerning

597     classification accuracy (McDermott et al., 2021).

598     While our study provides novel insight into the core functional components in decoding

599     multi-task fMRI signals, it is important to note that there are three limitations. The first

600     limitation is the manual setting of parameters for DBN and sparse representation framework,

31

601  mainly including the number of neuron nodes and layers in DBN and the sparsity penalty

602  parameter of SR. Thus, automatic optimization of model parameters is one of the future

603  research directions. The second limitation stems from our inability to detect FBNs related to

604  gambling and relational tasks within the first two to three layers of the DBN-SR framework.

605  This could be attributed to more noise present in the group-wise temporal features $\boldsymbol{D}^1$ extracted

606  at lower levels (Fig. 1). Additionally, LASSO regression may not be well-suited for handling

607  noisy shallow features, thus making it challenging for LASSO regression to accurately capture

608  the underlying spatial patterns. To address this limitation, future studies could explore

609  alternative regression approaches that are better suited for handling noisy shallow features,

610  thereby improving the accurate acquisition of the underlying spatial patterns. The third

611  limitation is that our study employed a relatively small dataset, consisting of 60 individuals out

612  of 68 from HCP Q1 dataset. To assess the robustness of our model, we included the remaining

613  8 individuals from the same dataset as a hold-out dataset, 6 of which do not have complete data

614  for all 7 tasks (sTab. 4). However, this does not affect their suitability as an independent lock

615  box dataset to test the performance of our trained model. The results revealed that the average

616  decoding accuracy for these 8 individuals (96.43%) was comparable to the 5-fold cross-

617  validation accuracy of the 60 individuals (sTab. 5), suggesting the robustness of our model.

618  Nonetheless, we acknowledge that a larger dataset would lend further support to our findings.

619  In future work, we aim to apply our model to more extensive or multicenter datasets to evaluate

620  its generalizability and robustness.

621  Overall, with the superiority of interpretability and effectiveness of DBN-SR model on

622  small datasets, our framework could potentially be useful to differentiate abnormal brain

623     function in clinical research.

634     **Reference**

635     Asja Fischer, C. I. (2012). An Introduction to Restricted Boltzmann Machines. Paper presented at the

636         Iberoamerican Congress on Pattern Recognition, Berlin.

637     Barch, D. M., Burgess, G. C., Harms, M. P., Petersen, S. E., Schlaggar, B. L., Corbetta, M., . . .

638         Consortium, W. U.-M. H. (2013). Function in the human connectome: task-fMRI and individual

639         differences in behavior. Neuroimage, 80, 169-189. doi:10.1016/j.neuroimage.2013.05.033

640     Ben J. Harrison, J. P., Marina Lo´ pez-Sola, Rosa Herna´ ndez-Ribas, Joan Deus, Hector Ortiz, Carles

641         Soriano-Mas, Murat Yu¨ cel, Christos Pantelis, and Narcı´s Cardoner. (2008). Consistency and

642    functional specialization in the default mode brain network. PNAS, 105, 9781–9786.

643    Benesty, J., Chen, J., Huang, Y., & Cohen, I. (2009). Pearson correlation coefficient. In Noise reduction

644        in speech processing (pp. 1-4): Springer.

645    Bengio, Y., Courville, A. C., & Vincent, P. (2012). Unsupervised feature learning and deep learning: A

646        review and new perspectives. CoRR, abs/1206.5538, 1(2665), 2012.

647    Bo Liu, Y. W., Yu Zhang, Qiang Yang. (2017, August). Deep Neural Networks for High Dimension, Low

648        Sample Size Data. Paper presented at the IJCAI, Melbourne.

649    Calhoun, V. D., Adali, T., McGinty, V. B., Pekar, J. J., Watson, T. D., & Pearlson, G. D. (2001). fMRI

650        activation in a visual-perception task: network of areas detected using the general linear model

651        and    independent    components    analysis.    Neuroimage,    14(5),    1080-1088.

652        doi:10.1006/nimg.2001.0921

653    Chang, C.-C., & Lin, C.-J. (2011a). Libsvm. ACM Transactions on Intelligent Systems and Technology,

654        2(3), 1-27. doi:10.1145/1961189.1961199

655    Chang, C.-C., & Lin, C.-J. (2011b). LIBSVM: a library for support vector machines. ACM transactions

656        on intelligent systems and technology (TIST), 2(3), 1-27.

657    Davatzikos, C., Ruparel, K., Fan, Y., Shen, D. G., Acharyya, M., Loughead, J. W., . . . Langleben, D. D.

658        (2005). Classifying spatial patterns of brain activity with machine learning methods: application

659        to lie detection. Neuroimage, 28(3), 663-668. doi:10.1016/j.neuroimage.2005.08.009

660    Dong, Q. (2020). Modeling Hierarchical Brain Networks via Volumetric Sparse Deep Belief Network

661        (VSDBN). Computerized Medical Imaging and Graphics.

662    Fangfei Ge, J. L., Xintao Hu , Lei Guo , Junwei Han , Shijie Zhao, Tianming Liu (2018, April 4-7).

663        Exploring intrinsic networks and their interactions using group wise temporal sparse coding.

664	Paper presented at the International Symposium on Biomedical Imaging (ISBI 2018),

665	Washington, D.C., USA.

666	Fisher, R. A., & Yates, F. (1938). Statistical tables for biological, agricultural aad medical research.

667	Statistical tables for biological, agricultural aad medical research.

668	Friston, K. J. (2009). Modalities, Modes, and Models in Functional Neuroimaging. SCIENCE, 326, 399-

669	403.

670	Hansen, K., Montavon, G., Biegler, F., Fazli, S., Rupp, M., Scheffler, M., . . . Muller, K. R. (2013).

671	Assessment and Validation of Machine Learning Methods for Predicting Molecular Atomization

672	Energies. J Chem Theory Comput, 9(8), 3404-3419. doi:10.1021/ct400195d

673	Haufe, S., Meinecke, F., Görgen, K., Dähne, S., Haynes, J.-D., Blankertz, B., & Bießmann, F. (2014).

674	On the interpretation of weight vectors of linear models in multivariate neuroimaging.

675	Neuroimage, 87, 96-110.

676	Haynes, J. D., & Rees, G. (2006). Decoding mental states from brain activity in humans. Nat Rev

677	Neurosci, 7(7), 523-534. doi:10.1038/nrn1931

678	Hinton, G. E., Osindero, S., & Teh, Y. W. (2006). A fast learning algorithm for deep belief nets. Neural

679	Comput, 18(7), 1527-1554. doi:10.1162/neco.2006.18.7.1527

680	Hinton, G. E., & Sejnowski, T. J. (1986). Learning and relearning in Boltzmann machines. Parallel

681	distributed processing: Explorations in the microstructure of cognition, 1(282-317), 2.

682	Hu, J., Kuang, Y., Liao, B., Cao, L., Dong, S., & Li, P. (2019). A Multichannel 2D Convolutional Neural

683	Network Model for Task-Evoked fMRI Data Classification. Comput Intell Neurosci, 2019,

684	5065214. doi:10.1155/2019/5065214

685	Hu, X., Huang, H., Peng, B., Han, J., Liu, N., Lv, J., . . . Liu, T. (2018). Latent source mining in FMRI

686     via restricted Boltzmann machine. Hum Brain Mapp, 39(6), 2368-2380. doi:10.1002/hbm.24005

687 Huan Liu , M. Z., Xintao Hu  , Yudan Ren  , Shu Zhang  , Junwei Han  , Lei Guo , Tianming Liu (2017).

688     Fmri data classification based on hybrid temporal and spatial sparse representation. Paper

689     presented at the IEEE 14th International Symposium on Biomedical Imaging (ISBI 2017),

690     Melbourne, VIC, Australia.

691 Huang, H., Hu, X., Zhao, Y., Makkie, M., Dong, Q., Zhao, S., . . . Liu, T. (2018). Modeling Task fMRI

692     Data Via Deep Convolutional Autoencoder. IEEE Trans Med Imaging, 37(7), 1551-1561.

693     doi:10.1109/TMI.2017.2715285

694 Huang, X., Xiao, J., & Wu, C. (2021). Design of Deep Learning Model for Task-Evoked fMRI Data

695     Classification. Comput Intell Neurosci, 2021, 6660866. doi:10.1155/2021/6660866

696 Jang, H., Plis, S. M., Calhoun, V. D., & Lee, J. H. (2017). Task-specific feature extraction and

697     classification of fMRI volumes using a deep neural network initialized with a deep belief network:

698     Evaluation using sensorimotor tasks. Neuroimage, 145(Pt B), 314-328.

699     doi:10.1016/j.neuroimage.2016.04.003

700 Kay, K., Rokem, A., Winawer, J., Dougherty, R., & Wandell, B. (2013). GLMdenoise: a fast, automated

701     technique for denoising task-based fMRI data. Frontiers in neuroscience, 247.

702 Kohavi, R. (1995). A study of cross-validation and bootstrap for accuracy estimation and model

703     selection. Paper presented at the Ijcai.

704 Kriegeskorte, N., & Bandettini, P. (2007). Analyzing for information, not activation, to exploit high-

705     resolution fMRI. Neuroimage, 38(4), 649-662.

706 LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. Nature, 521(7553), 436-444.

707     doi:10.1038/nature14539

708    Lee, J., Jeong, Y., & Ye, J. C. (2013). Group sparse dictionary learning and inference for resting-state

709        fMRI analysis of Alzheimer's disease. Paper presented at the 2013 IEEE 10th International

710        Symposium on Biomedical Imaging.

711    Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., . . . Sanchez, C. I.

712        (2017). A survey on deep learning in medical image analysis. Med Image Anal, 42, 60-88.

713        doi:10.1016/j.media.2017.07.005

714    Liu, X., He, P., Chen, W., & Gao, J. (2019). Multi-task deep neural networks for natural language

715        understanding. arXiv preprint arXiv:1901.11504.

716    Logothetis, N. K. (2008). What we can do and what we cannot do with fMRI. Nature, 453(7197), 869-

717        878.

718    Lv, J., Jiang, X., Li, X., Zhu, D., Chen, H., Zhang, T., . . . Liu, T. (2015). Sparse representation of whole-

719        brain fMRI signals for identification of functional networks. Med Image Anal, 20(1), 112-134.

720        doi:10.1016/j.media.2014.10.011

721    Lv, J., Jiang, X., Li, X., Zhu, D., Zhang, S., Zhao, S., . . . Liu, T. (2015). Holistic atlases of functional

722        networks and interactions reveal reciprocal organizational architecture of cortical function. IEEE

723        Trans Biomed Eng, 62(4), 1120-1131. doi:10.1109/TBME.2014.2369495

724    McDermott, E. J., Raggam, P., Kirsch, S., Belardinelli, P., Ziemann, U., & Zrenner, C. (2021). Artifacts

725        in EEG-Based BCI Therapies: Friend or Foe? Sensors (Basel), 22(1). doi:10.3390/s22010096

726    Najafabadi, M. M., Villanustre, F., Khoshgoftaar, T. M., Seliya, N., Wald, R., & Muharemagic, E. (2015).

727        Deep learning applications and challenges in big data analytics. Journal of big data, 2(1), 1-21.

728    O'Reilly, J. X., Woolrich, M. W., Behrens, T. E., Smith, S. M., & Johansen-Berg, H. (2012). Tools of the

729        trade: psychophysiological interactions and functional connectivity. Social cognitive and

730       affective neuroscience, 7(5), 604-609.

731    Qiang, N., Dong, Q., Zhang, W., Ge, B., Ge, F., Liang, H., . . . Liu, T. (2020). Modeling task-based fMRI

732       data via deep belief network with neural architecture search. Comput Med Imaging Graph, 83,

733       101747. doi:10.1016/j.compmedimag.2020.101747

734    Rashid, M., Singh, H., & Goyal, V. (2020). The use of machine learning and deep learning algorithms

735       in functional magnetic resonance imaging—a systematic review. Expert Systems, 37(6),

736       e12644. doi:10-1111

737    Ren, Y., Xu, S., Tao, Z., Song, L., & He, X. (2021). Hierarchical Spatio-Temporal Modeling of

738       Naturalistic Functional Magnetic Resonance Imaging Signals via Two-Stage Deep Belief

739       Network With Neural Architecture Search. Front Neurosci, 15, 794955.

740       doi:10.3389/fnins.2021.794955

741    Rubin, T. N., Koyejo, O., Gorgolewski, K. J., Jones, M. N., Poldrack, R. A., & Yarkoni, T. (2017).

742       Decoding brain activity using a large-scale probabilistic functional-anatomical atlas of human

743       cognition. PLoS Comput Biol, 13(10), e1005649. doi:10.1371/journal.pcbi.1005649

744    Salimi-Khorshidi, G., Douaud, G., Beckmann, C. F., Glasser, M. F., Griffanti, L., & Smith, S. M. (2014).

745       Automatic denoising of functional MRI data: combining independent component analysis and

746       hierarchical fusion of classifiers. Neuroimage, 90, 449-468.

747       doi:10.1016/j.neuroimage.2013.11.046

748    Shu Zhang , X. L., Lei Guo , Tianming Liu. (2017, 18-21 April). Exploring human brain activation via

749       nested sparse coding and functional operators. Paper presented at the International

750       Symposium on Biomedical Imaging (ISBI 2017), Melbourne, VIC, Australia.

751    Song, L., Ren, Y., Hou, Y., He, X., & Liu, H. (2022). Multitask fMRI Data Classification via Group-Wise

752        Hybrid Temporal and Spatial Sparse Representations. eNeuro, 9(3).

753        doi:10.1523/ENEURO.0478-21.2022

754    Sotetsu Koyamadaa, b., Yumi Shikauchia,b, Ken Nakaea, Masanori Koyamaa, Shin Ishii. (2015). Deep

755        learning of fMRI big data: a novel approach to subject-transfer decoding. arXiv preprint arXiv.

756    Stanislas Dehaene, G. L. C. H., Laurent Cohen, Jean-Baptiste Poline, Pierre-François van de Moortele

757        and Denis Le Bihan. (1998). Inferring behavior from functional brain images.

758    Tibshirani, R. ( 2011). Regression shrinkage and selection via the lasso:

759    a retrospective. Royal Statistical Society, 73, 273-282.

760    Varoquaux, G., & Thirion, B. (2014). How machine learning is shaping cognitive neuroimaging.

761        GigaScience, 3(1), 1-7. doi:10.1186

762    Vieira, S., Pinaya, W. H., & Mechelli, A. (2017). Using deep learning to investigate the neuroimaging

763        correlates of psychiatric and neurological disorders: Methods and applications. Neurosci

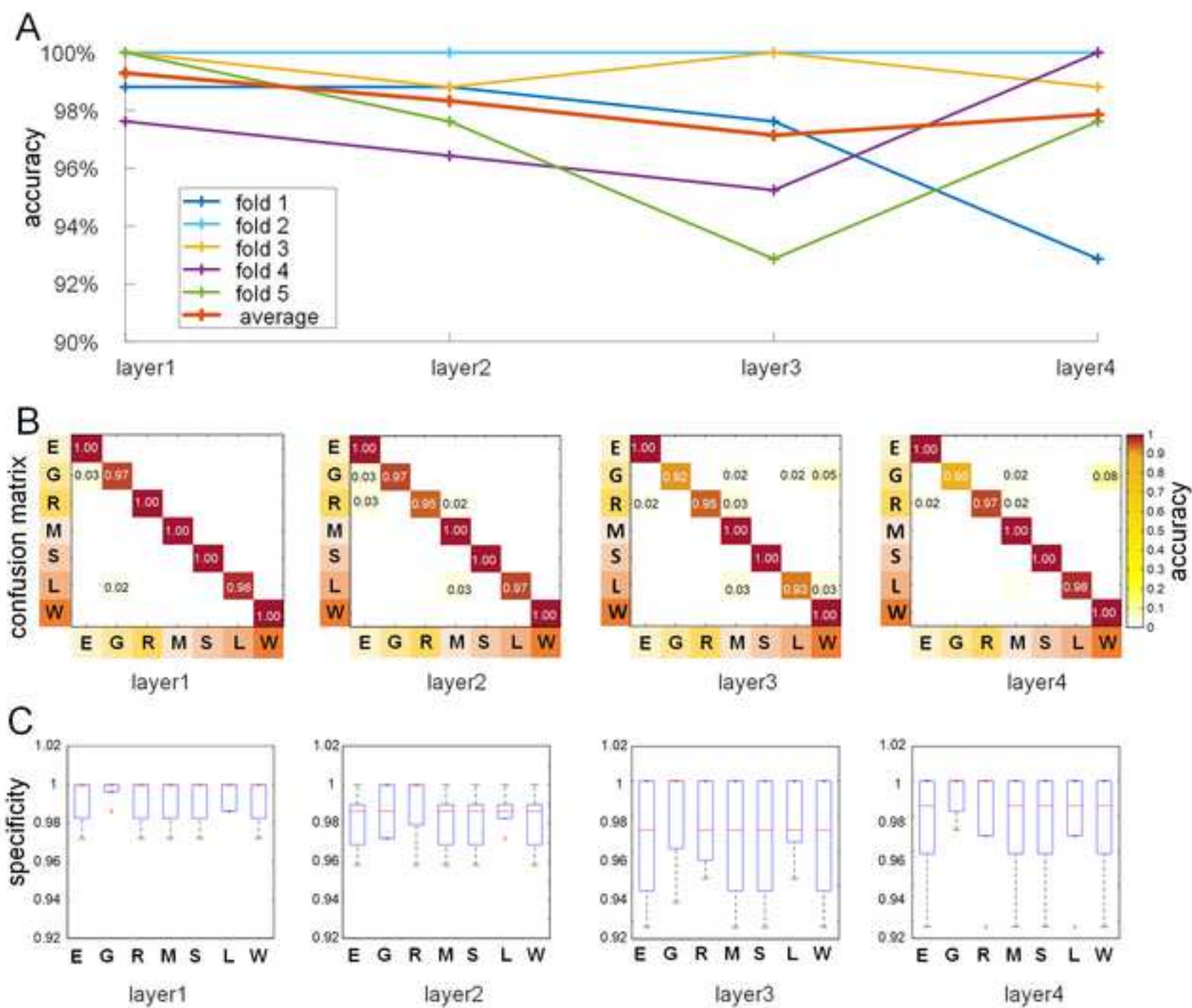764        Biobehav Rev, 74(Pt A), 58-75. doi:10.1016/j.neubiorev.2017.01.002

765    Wang, X., Liang, X., Jiang, Z., Nguchu, B. A., Zhou, Y., Wang, Y., . . . Qiu, B. (2020). Decoding and

766        mapping task states of the human brain via deep learning. Hum Brain Mapp, 41(6), 1505-1519.

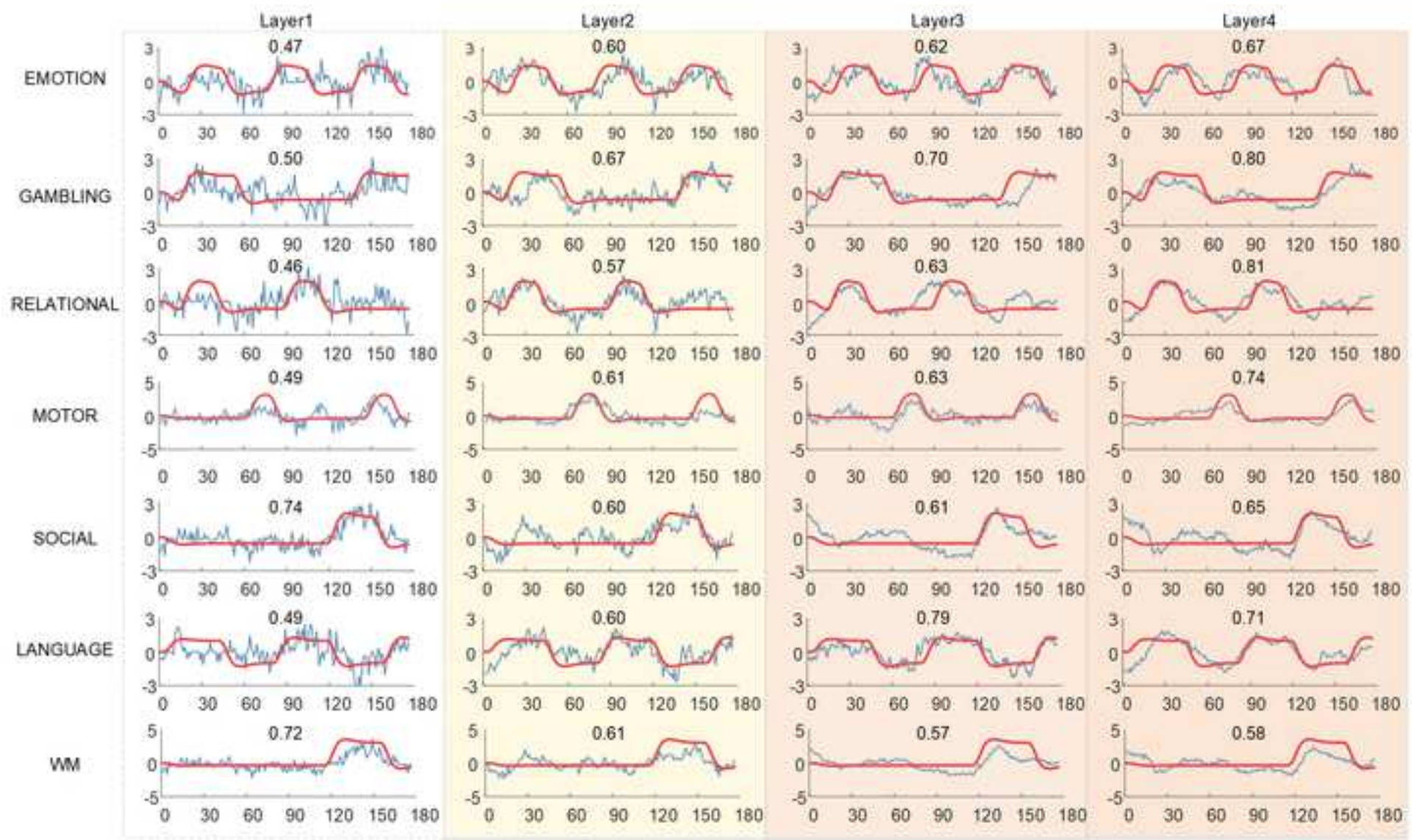767        doi:10.1002/hbm.24891

768    Wei Zhang, S. Z., Xintao Hu,2, Qinglin Dong,Heng Huang,Shu Zhang, Yu Zhao, Haixing Dai, Fangfei

769        Ge, Lei Guo and Tianming Liu. (2020). Hierarchical Organization of Functional Brain Networks

770        Revealed by Hybrid Spatiotemporal Deep Learning. Brain Connectivity, 10.

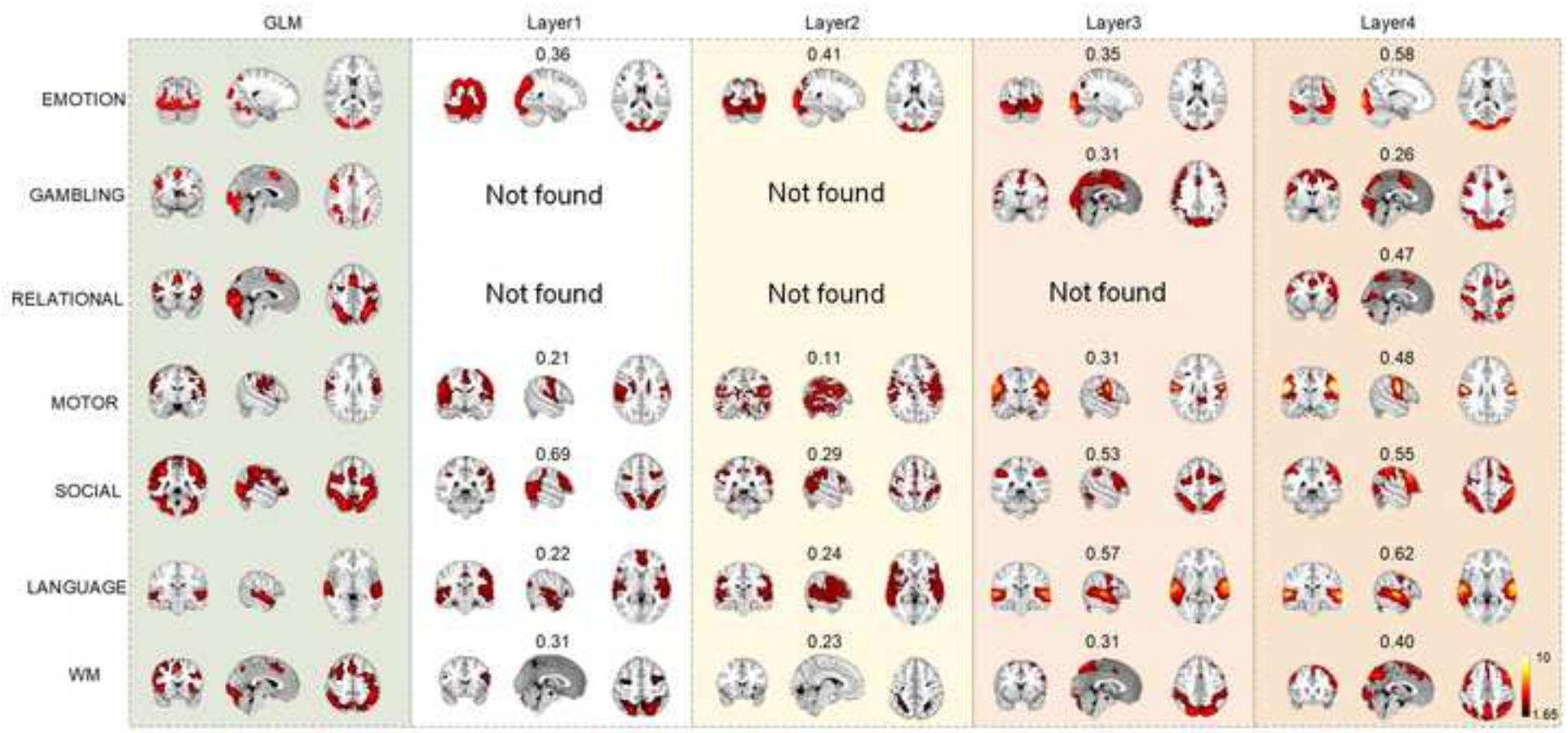771        doi:10.1089/brain.2019.0701

772    Wen, D., Wei, Z., Zhou, Y., Li, G., Zhang, X., & Han, W. (2018). Deep Learning Methods to Process

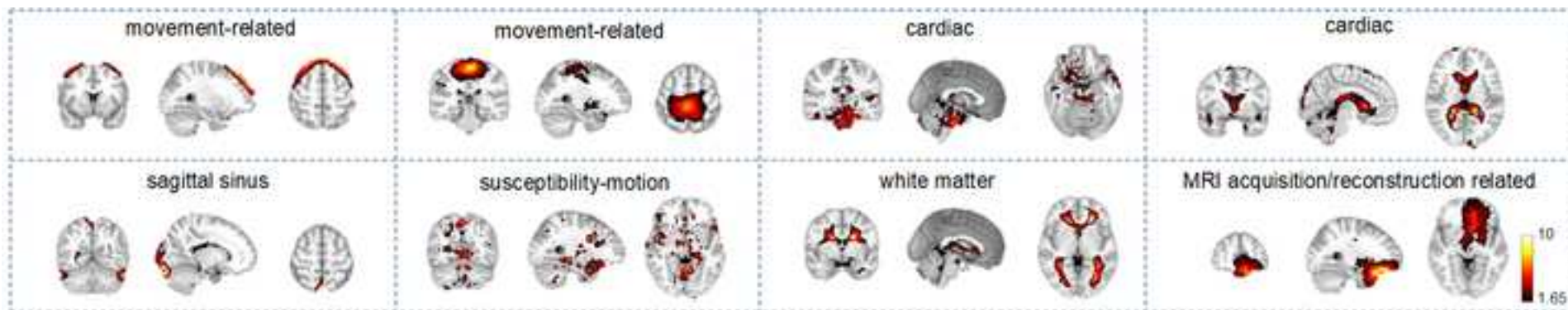773        fMRI Data and Their Application in the Diagnosis of Cognitive Impairment: A Brief Overview

774      and Our Opinion. Front Neuroinform, 12, 23. doi:10.3389/fninf.2018.00023

775    Xu, S., Ren, Y., Tao, Z., Song, L., & He, X. (2022). Hierarchical Individual Naturalistic Functional Brain

776      Networks with Group Consistency uncovered by a Two-Stage NAS-Volumetric Sparse DBN

777      Framework. eNeuro, 9(5). doi:10.1523/ENEURO.0200-22.2022

778    Zhang, S., Li, X., Lv, J., Jiang, X., Guo, L., & Liu, T. (2016). Characterizing and differentiating task-

779      based and resting state fMRI signals via two-stage sparse representations. Brain Imaging

780      Behav, 10(1), 21-32. doi:10.1007/s11682-015-9359-7

781    Zhang, Y., Tetrel, L., Thirion, B., & Bellec, P. (2021). Functional annotation of human cognitive states

782      using deep graph convolution. Neuroimage, 231, 117847.

783      doi:10.1016/j.neuroimage.2021.117847

784    Zuo, X. N., Kelly, C., Adelstein, J. S., Klein, D. F., Castellanos, F. X., & Milham, M. P. (2010). Reliable

785      intrinsic connectivity networks: test-retest evaluation using ICA and dual regression approach.

786      Neuroimage, 49(3), 2163-2177. doi:10.1016/j.neuroimage.2009.10.080
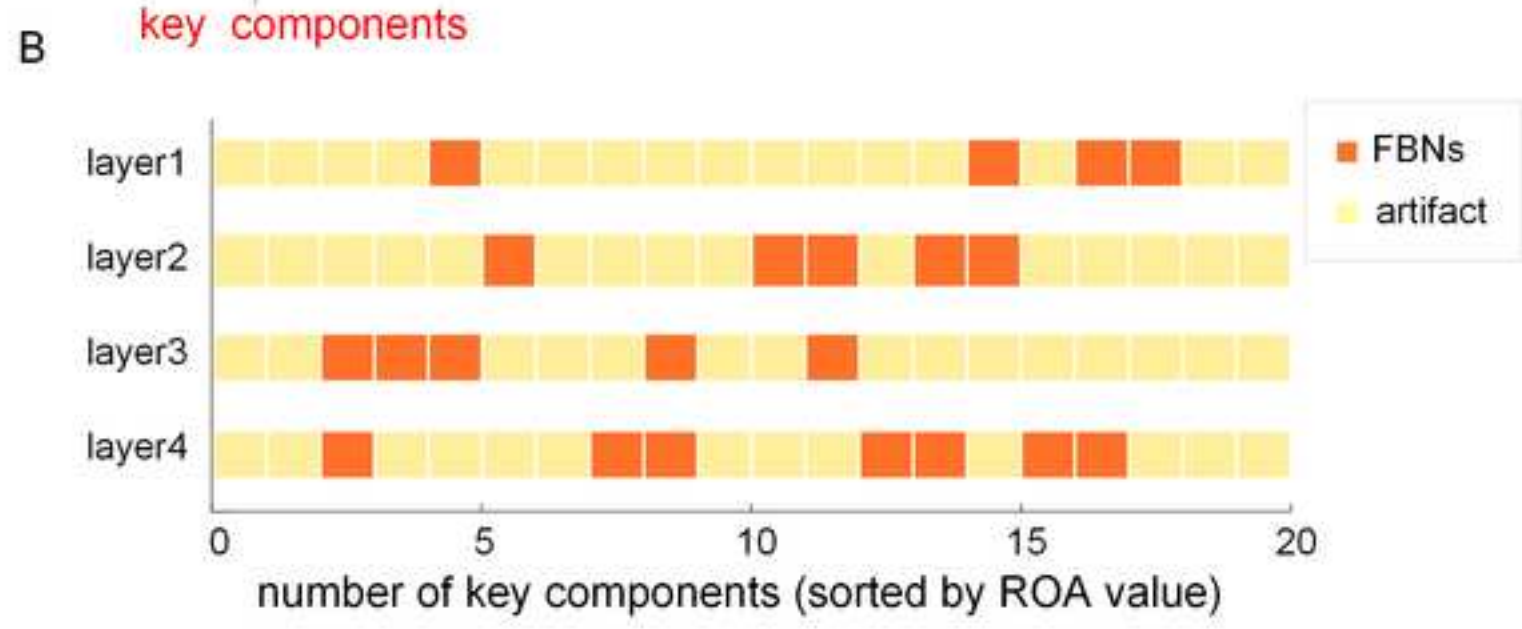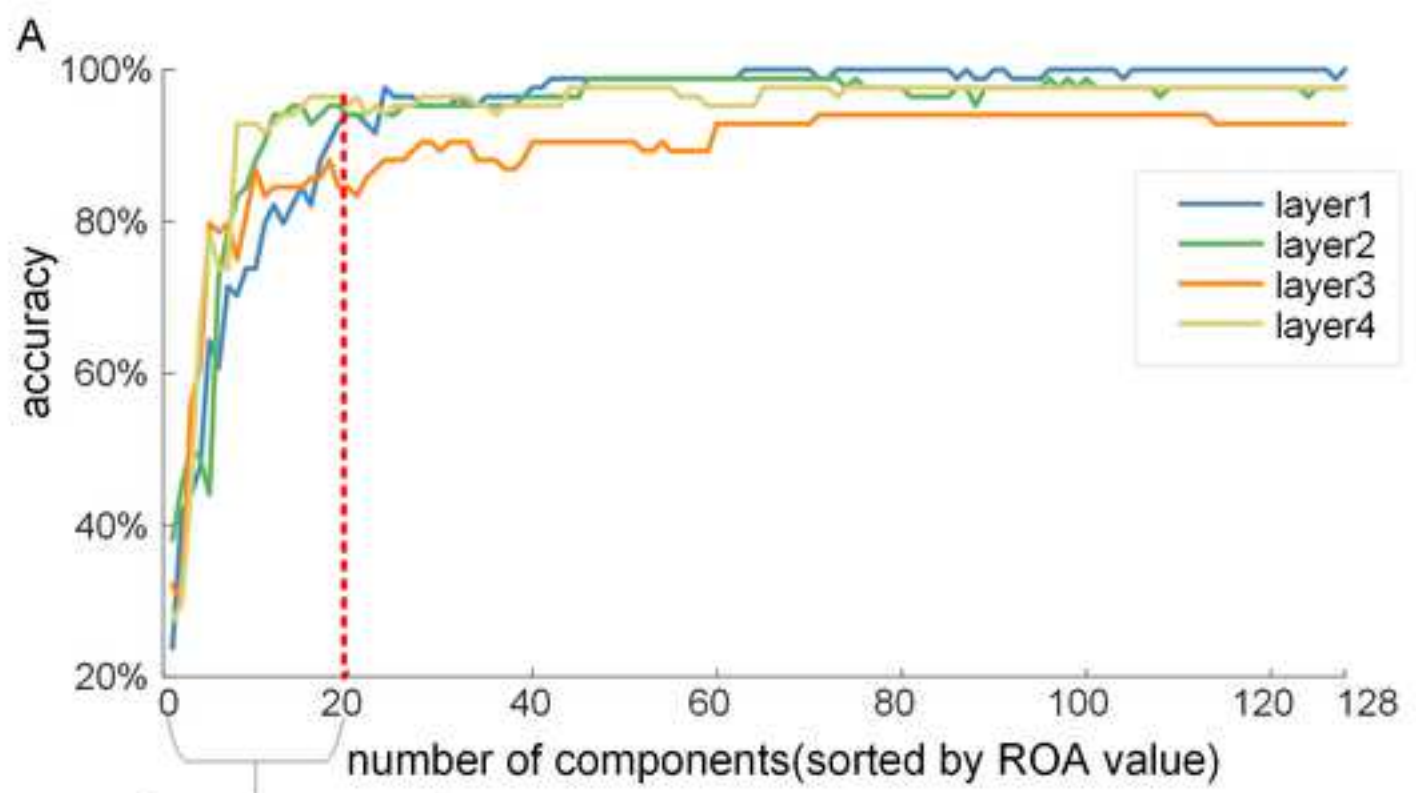
# Author Summary

Decoding different cognitive processes using task-based functional magnetic resonance imaging (tfMRI) is crucial for understanding the relationship between brain activities and cognitive states. However, existing machine learning-based feature extraction methods for decoding brain states may struggle to capture the complex and precise spatiotemporal patterns of brain activity from the highly noisy raw fMRI data. Additionally, current deep learning-based end-to-end decoding models struggle to unveil interpretable components in tfMRI signal decoding.

To address these limitations, we proposed a novel framework, the hybrid spatio-temporal deep belief network and sparse representations (DBN-SR) framework, which effectively distinguished multi-task fMRI signals with an average accuracy of 97.86%. Furthermore, it simultaneously identified multi-level temporal and spatial patterns of multiple cognitive tasks. By utilizing a novel Ratio-of-Activation metric, our framework unveiled interpretable components with greater classification capacity, offering an effective methodology for basic neuroscience and clinical research.