# Deep Linear Modeling of Hierarchical Functional Connectivity in the Human Brain

Wei Zhang[1], Eva Palacios[1], Pratik Mukherjee[1,2#]

[1]Department of Radiology and Biomedical Imaging,

[2]Department of Bioengineering and Therapeutic Sciences,

University of California San Francisco,

San Francisco, CA, 94143-0628, USA;

E-mail: pratik.mukherjee@ucsf.edu

[#]Corresponding author

**Abstract**

The human brain exhibits hierarchical modular organization, which is not depicted by conventional fMRI functional connectivity reconstruction methods such as independent component analysis (ICA). To map hierarchical brain connectivity networks (BCNs), we propose a novel class of deep (multilayer) linear models that are constructed such that each successive layer decomposes the features of the preceding layer. Three of these are multilayer variants of Sparse Dictionary Learning (SDL), Non-Negative Matrix Factorization (NMF) and Fast ICA (FICA). We present a fourth deep linear model, Deep Matrix Fitting (MF), which incorporates both rank reduction for data-driven hyperparameter determination as well as a distributed optimization function. We also introduce a novel framework for theoretical comparison of these deep linear models based on their combination of mathematical operators, the predictions of which are tested using simulated resting state fMRI data with known ground truth BCNs. Consistent with the theoretical predictions, Deep MF and Deep SDL performed best for connectivity estimation of $1^{st}$ layer networks, whereas Deep FICA and Deep NMF were modestly better for spatial mapping. Deep MF provided the best overall performance, including computational speed. These deep linear models can efficiently map hierarchical BCNs without requiring the manual hyperparameter tuning, extensive fMRI training data or high-performance computing infrastructure needed by deep nonlinear models, such as convolutional neural networks (CNNs) or deep belief networks (DBNs), and their results are also more explainable from their mathematical structure. These benefits gain in importance as continual improvements in the spatial and temporal resolution of fMRI reveal more of the hierarchy of spatiotemporal brain architecture. These new models of hierarchical BCNs may also advance the development of fMRI diagnostic and prognostic biomarkers, given the recent recognition of disparities between low-level vs high-level network connectivity across a wide range of neurological and psychiatric disorders.

1                                 Abbreviations

2

3    ADMM: Alternating Direction Method of Multipliers

4    BCN: Brain Connectivity Network

5    CNN: Convolutional Neural Network

6    CPU: Central Processing Unit

7    DBN: Deep Belief Network

8    DCAE: Deep Convolutional Auto Encoder

9    Deep FICA: Deep Fast Independent Component Analysis

10    Deep MF: Deep Matrix Fitting

11    Deep NMF: Deep Non-negative Matrix Factorization

12    Deep SDL: Deep Sparse Dictionary Learning

13    DNN: Deep Neural Network

14    fMRI: Functional Magnetic Resonance Imaging

15    GD: Gradient Descent

16    GLM: General Linear Model

17    GPU: Graphics Processing Unit

18    HD: Hausdorff Distance

19    ICA: Independent Component Analysis

20    IS: Intensity Similarity

21    LASSO: Least Absolute Shrinkage and Selection Operator

22    RBM: Restricted Boltzmann Machine

23    RRO: Rank Reduction Operator

24    rsfMRI: Resting-State Functional MRI

25    SS: Spatial Similarity

26    tfMRI: Task-Evoked Functional MRI

27    TBI: Traumatic Brain Injury

28    TPU: Tensor Processing Unit

29

## 1. Introduction

Functional Magnetic Resonance Imaging (fMRI) has been widely used for the identification of brain connectivity networks (BCNs) (Bartels et al., 2005; Beckmann et al., 2005; Biswal et al., 1995, 2010; Bullmore et al., 2009; Duncan et al., 2010; Stam et al., 2014). A variety of scientific investigations have already demonstrated the hierarchical modular organization of human brain networks (Bassett et al., 2008; Biswal et al., 2005; Bullmore et al., 2009; Sporns et al., 2004). The architecture of cortical BCNs is organized at different spatial scales, from both functional and structural perspectives, ranging from local circuits at the microscale to columns and layers at the mesoscale to areas and areal networks at the macroscale (Bullmore et al., 2009; Power et al., 2011; Stam et al., 2014; Sporns et al., 2004).

In the last two decades, a variety of computational methods have been developed to detect BCNs, e.g., General Linear Modeling (GLM), Graph Theory, Independent Component Analysis (ICA), and Sparse Dictionary Learning (SDL) (Andersen et al., 1999; Calhoun et al., 2001; Lee et al., 2011; Lee et al., 2016; Lv et al., 2015; Zhang et al., 2017; Zhang et al., 2018; Zhang et al., 2019). However, these methods are based on a 'shallow' framework that cannot identify in unsupervised data-driven fashion the hierarchical and spatially overlapping organization of BCNs using resting-state fMRI (rsfMRI) or task-evoked fMRI (tfMRI) signals (Hu et al., 2018; Huang et al., 2018; Zhang et al., 2019; Zhang et al., 2020). Traditionally, the hierarchical spatial organization of BCNs has been indicated by varying the number of features in shallow linear models, for example from low to high numbers of independent components in ICA (Iraji et al., 2019; Smith et al., 2009), and noting that smaller networks at the more granular decomposition tend to merge or otherwise recombine to form larger networks at the coarser decomposition. However, there is no principled, unsupervised way to map this hierarchical organization with shallow methods.

Fortunately, with the advent of deep learning, algorithms have been developed that are capable of reconstructing hierarchical network architectures, e.g., the Deep Convolutional Auto Encoder (DCAE), Deep Belief Network (DBN) and Convolutional Neural Network (CNN)

4

1  (Bengio et al., 2012; Esteva, et al., 2019; Gurovich et al., 2019; Hannun et al., 2019; LeCun

2  et al., 2015; Plis et al., 2014; Schmidhuber et al., 2015; Suk et al., 2014; Suk et al., 2016;

3  Zhang et al., 2020). The Restricted Boltzmann Machine (RBM) can be used to model fMRI

4  time series signals and effectively reconstruct functional brain networks with impressive

5  accuracy (Hu et al., 2018; Huang et al., 2018). Moreover, other recent studies reported

6  meaningfully hierarchical temporal organization of tfMRI time series, each with corresponding

7  task-evoked BCNs (Hu et al., 2018; Zhang et al., 2019; Zhang et al., 2020) using DCAE, RBM

8  and DBN. In general, these machine learning techniques are considered to be deep nonlinear

9  models, e.g., deep neural networks (DNN). Although these nonlinear models such as DBN

10  have recently proven effective at hierarchical spatiotemporal decomposition of task-evoked

11  fMRI data (Dong et al., 2020), there are several disadvantages: (i) large training samples; (ii)

12  extensive computational resources, e.g., graphics processing units (GPUs) or tensor

13  processing units (TPUs); (iii) manual tuning of hyperparameters; (iv) time-consuming training

14  process; (v) non-convergence to the global optimum; and (vi) "black box" results that lack

15  explainability. Deep linear algorithms can overcome all these shortcomings of nonlinear

16  techniques, since they are fast even on conventional central processing units (CPUs) with

17  hyperparameters that can be automatically determined and with convex optimization functions

18  that are guaranteed to converge. Furthermore, as we show in the theoretical analysis below,

19  important aspects of their behavior can be explained from their relatively simple mathematical

20  structure. For fMRI research, these deep linear models can detect BCNs using data from

21  relatively few experimental subjects compared to deep nonlinear models and may prove

22  especially useful and efficient as the spatial and temporal resolution of fMRI continues to

23  improve, revealing more of the hierarchy of brain organization.

24

25      For hierarchical spatial functional connectivity mapping, we adopt a compositional

26  approach to develop multilayer versions of SDL (Deep SDL), Fast ICA (Deep FICA; Seo, 2018)

27  and Non-negative Matrix Factorization (Deep NMF; Trigeorgis et al., 2016), as well as a novel

28  multilayer linear model that we name Deep Matrix Fitting (Deep MF). We contrast these deep

29  linear algorithms for fMRI functional connectivity analysis in two ways. First, we employ theory

30  to investigate the mathematical properties of these models, in order to predict differences in,

1    e.g., network sparsity, connectivity strength and convergence velocity. Second, we conduct *in*

2    *silico* connectivity reconstruction experiments using simulated fMRI signal time series to test

3    the predictions of the theoretical analyses for the relative performance of the four deep linear

4    models for fMRI brain network mapping. This leads to clear conclusions about the strengths

5    and weaknesses of each method and provides a guide to the research community for applying

6    this novel class of network reconstruction methodologies as well as for developing new ones.
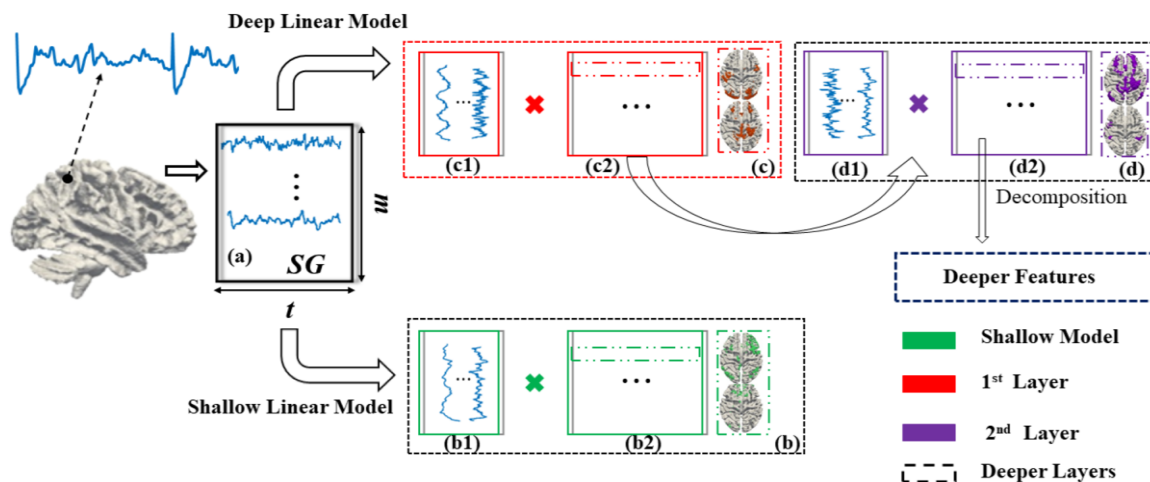
7

8    ## 2.    Methods

9    *2.1 Shallow versus Deep Linear Models of fMRI Functional Connectivity*

10    The following introductions in Sections 2.2 to 2.5 provide the fundamentals of each deep

11    linear model. Furthermore, these descriptions are prerequisites to analyze the theoretical

12    properties of each model in the succeeding sections. Figure 1 compares the computational

13    steps of a deep linear model versus that of a conventional "shallow" linear model:

14



15    **Figure 1.** Deep linear model versus shallow linear model. The shallow model has only a single layer,
16    i.e., a single decomposition. The deep linear model is constructed via multiple layers, i.e., continuous
17    decomposition. (**a**) $SG$ represents the input fMRI signal matrix; it contains the $t$ time points and $m$ voxels.
18    (**b**) describes the pipeline of a shallow model in which the original input signal is decomposed into the
19    weight matrix/dictionary (shown as **b1**) and feature matrix, i.e. connectivity networks (shown as **b2**). (**c**)
20    and (**d**) represent the 1st and 2nd layers of the linear deep model, respectively. (**c1**) represents the 1st
21    layer weight matrix/dictionary identified via SG. (**c2**) represents the 1st layer feature matrix, i.e.
22    connectivity networks, recognized via SG. Similarly, (**d1**) and (**d2**) represent the corresponding matrices
23    of the 2nd layer, which are both derived from the 1st layer feature matrix. The dashed blue rectangle
24    indicates the deeper features beyond the 2nd layer that are derived from the 2nd layer feature matrix.

25

1     *2.2 Deep Matrix Fitting*

2        We propose a novel and efficient deep linear model that we name Deep Matrix Fitting. This

3     algorithm aims to detect the hierarchical and overlapping organization of BCNs better than

4     previously described data-driven functional connectivity reconstruction methods, e.g., ICA,

5     SDL, DCAE and DBN (Calhoun et al., 2001; Lv et al., 2015; Zhang et al., 2020; Hinton and

6     Salakhutdinov, 2006; Hinton et al., 2012). Due to the constraints of spatial independence in

7     ICA, some investigators have reported that ICA cannot easily identify extensively overlapped

8     functional brain networks (Calhoun et al., 2001; McKeown and Sejnowski, 1998; Zhang et al.,

9     2019). Although SDL can efficiently derive spatial features, i.e., functional brain networks,

10     based on rsfMRI and tfMRI, it is very challenging to leverage the dictionary size, sparsity trade

11     off and even number of layers to implement a deep SDL. To be specific, one must heuristically

12     estimate the dictionary size and number of layers. Simply utilizing the same size of dictionary

13     and number of layers can easily result in the vanishing of spatial features of deeper layers,

14     due to iteratively using the $\ell_1$ norm. Recent deep nonlinear models, such as DBN, can

15     successfully reveal the architecture of hierarchical spatiotemporal features. Unfortunately, the

16     probabilistic energy-based model of DBNs necessarily requires a large number of training

17     samples to avoid overfitting. Furthermore, DBN requires extensive computational resources

18     such as GPUs and even TPUs (Zhang et al., 2019; Zhang et al., 2020). The novel Deep MF

19     proposed in this work successfully solves these aforementioned problems. Deep MF can

20     automatically estimate the optimal dictionary size, sparsity trade-off and number of layers,

21     using an operator of rank reduction (Wen et al., 2012; Shen et al., 2014). In other words, Deep

22     MF does not require any manual hyperparameter tuning to decompose the rsfMRI signal

23     matrix. Since Deep MF is a deep linear model, it should detect latent features faster than DBN

24     while only requiring conventional CPUs. In general, Deep MF can be approximately

25     considered as a deep SDL (described in Section 2.3) with the additional mechanism to

26     automatically determine all crucial hyperparameters via rank reduction.

27

28        The equation governing Deep MF is:

7

$$min_{X_i, \quad Y_i, \quad S \in \mathbb{R}^{m \times n}} \left\| \prod_{i=1}^{M-1} X_i Y_M - SG \right\|_F^2 + \mu \sum_{i=1}^{M} \|Y_i\|_1 + \lambda \sum_{i=1}^{M} \|Z_i\|_1 \qquad (1)$$

$$X_i Y_i \leftarrow \mathcal{R}(Y_{i-1})$$

1   where $\{X_i\}_{i=1}^{M}$ represents the hierarchical dictionaries, e.g., $X_i$ indicates the dictionary of the

2   $i$ th layer. $\{X_i\}_{i=1}^{M}$ is also considered as the time series in GLM and the weight matrix in ICA

3   and DBN. $M$ is the total number of layers. Similarly, $\{Y_i\}_{i=1}^{M}$ represents the hierarchical

4   spatial features, e.g., $Y_i$ indicates the spatial features of $i^{th}$ layer. $\{Y_i\}_{i=1}^{M}$ is also denoted as a

5   correlation matrix. $\{Z_i\}_{i=1}^{M}$ are the matrices of background components, which is usually

6   treated as the noise. $\mathcal{R}$ represents a rank reduction operator (RRO) to automatically estimate

7   the hyperparameters and more details will be introduced in the following section. Naturally, we

8   assume the spatial features $Y_{i-1}$ can be decomposed as deeper dictionary $X_i$ and spatial

9   features $Y_i$, in order to implement the deep linear framework (Figure 1). Therefore, the original

10  input data $SG$ can be decomposed as $\prod_{i=1}^{M-1} X_i Y_M$. In Eq. (1), $\lambda$ and $\mu$ are known as the

11  sparse trade off to control the sparsity levels of background components and spatial features,

12  respectively. In addition, in Eq. (1), $\|\cdot\|_F$ and $\|\cdot\|_1$ represent the Frobenius and $\ell_1$ norms,

13  respectively.

14

15      This optimization function, shown as Eq. (1), consists of more parameters than ICA and

16  SDL. In general, SDL includes two parameters to be optimized: dictionary and correlation

17  matrix. Naturally, it is easier to comprehensively employ alternative optimizer and shrinkage

18  methods (Wen et al., 2012). Before optimizing Eq. (1), we need to convert Eq. (1) to an

19  augmented Lagrangian function. If considering the $k^{th}$ layer, we have:

$$\mathcal{L}_\beta \left( \prod_{i=1}^{k-1} X_i, Y_k, Z_k, e_k \right) \overset{\text{def}}{=} \left\| \prod_{i=1}^{k-1} X_i Y_k - SG \right\|_F^2 + \left\langle \prod_{i=1}^{k-1} X_i Y_k - SG, e_k \right\rangle \qquad (2)$$

20      For Eq. (2), for $k$ layers (we assume the total number of layers as $k$), these can be solved

21  using Alternating Direction of Method of Multipliers (ADMM) (Shen et al., 2014), and to solve

22  $\sum_{i=1}^{k} \|Z_i\|_1$, we jointly utilize the shrinkage method. In Eq. (2), all parameters are as discussed

23  before, with $e_k$ defined as the multiplier. The $\ell_1$ norm of $Y_k$ and $Z_k$ shown in Eq. (1) can be

24  solved directly using the shrinkage method (Beck et al., 2009).

25

1    The iterative format to solve Eq. (2) using ADMM can be organized as follows:

$$X_k^{it+1} = argmin_{X_k^{it+1} \in \mathbb{R}^{m \times h_k}} \; \mathcal{L}_\beta(X_k^{it}, Y_k^{it}, Z_k^{it}, e_k^{it}) \tag{3-1}$$

$$Y_k^{it+1} = argmin_{Y_k^{it+1} \in \mathbb{R}^{h_k \times n}} \; \mathcal{L}_\beta(X_k^{it+1}, Y_k^{it}, Z_k^{it}, e_k^{it}) \tag{3-2}$$

$$Z_k^{it+1} = argmin_{Z_k^{it+1} \in \mathbb{R}^{m \times n}} \; \mathcal{L}_\beta(X_k^{it+1}, Y_k^{it+1}, Z_k^{it}, e_k^{it}) \tag{3-3}$$

$$e_k^{it+1} = \; e_k^{it} + \beta(\prod_{i=1}^{k-1} X_i Y_k + \sum_{i=1}^{k} Z_k^{it+1} - SG) \tag{3-4}$$

2    Using ADMM, in each iteration (the current iteration is represented as *it*), we update a

3    single parameter independently, and finally calculate the multiplier, based on the current error.

4    Since, in each single step, only one parameter is optimized and others are fixed, Eq (2) is

5    considered as a convex problem and the global optimum can be obtained via a descent

6    algorithm, e.g. gradient descent (GD) or ADMM. In Eq. (3-4), $\beta$ denotes the step length.

7

8    To automatically estimate the dictionary size and number of layers, we introduce the

9    operator RRO. Briefly, RRO focuses on the identification of major components included in the

10    raw data, and simultaneously determines which components are relatively weak and that will

11    therefore be continuously merged into the background matrices. In general, RRO

12    demonstrates that the number of units, i.e., dictionary size, should be consistently reduced, if

13    considering deeper layers (Hinton et al., 2012; Zhang et al., 2019). In other words, the

14    continuous increase of units in deeper layers can result in lack of convergence. If the number

15    of units/dictionary size, i.e., the estimated rank of the matrix, is reduced to one, that indicates

16    the decomposition should be terminated. Hence, the layer that owns a rank of unity should be

17    considered the final layer. Deep MF employs RRO to continuously reduce the dictionary size

18    and therefore also determine the number of layers. In fact, it does not require any manual

19    design for the essential hyperparameters of deep learning models, such as the number of

20    layers or unit number of each layer that are used in DBN and other peer deep models.

21

22    In detail, this rank estimator RRO employs a technique of rank-revealing by continuously

23    using orthogonal decomposition, in this case via *QR* factorization (Wen et al., 2012; Shen et

24    al., 2014). The advantage of *QR* is that it is faster and makes fewer requirements of the input

1   matrix. For example, *QR* performs orthogonal decomposition faster than Singular Value

2   Decomposition (*SVD*) and can solve incomplete (i.e., number of features < number of samples)

3   and over-complete (i.e., number of features > number of samples) matrices.

4

5   At the beginning, $r^*$ is denoted as the initial estimated rank of $S^i$ and we denote *r* as the

6   optimal rank estimation of input matrix $S^i$. If $r^* \geq r$ holds, the detection of the diagonal line of

7   the upper-triangular matrix in the *QR* factorization can be performed using the input matrix $S^i$.

8   If we can determine the ideal size of *QR* factorization using $S^i$ in the work with permutation

9   matrix *E*, the diagonal matrix *R* is non-increasing in magnitude (Wen et al., 2012; Shen et al.,

10   2014). The *QR* factorization and rank-revealing will eventually provide a reasonable solution

11   using a proper thresholding value introduced in Eq. (2) and Eq. (3) (Wen et al., 2012; Shen et

12   al., 2014). By detecting the diagonal line of matrix *R*, we compute two vectors $d \in \mathbb{R}^r$ and

13   $r \in \mathbb{R}^{r-1}$:

$$d_i = |R_{ii}|$$
$$r_i = \frac{d_i}{d_{i+1}} \tag{4}$$

14   And then examine the value:

$$\xi = \frac{(m-1)r(p)}{\sum_{i \neq p} r_i} \tag{5}$$

15   where *r(p)* is the maximum element of the vector *r* (with the largest index *p* if the maximum

16   value is not unique). In our current implementation, we reset the rank estimated *r*, if $\xi > 2$,

17   and this adjustment can be successfully done only once (Wen et al., 2012; Shen et al., 2014).

18

19   The mathematical definition of RRO is shown below:

$$\mathcal{R}\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_{n-2}^{(1)} \\ a_{n-1}^{(1)} \end{bmatrix} \quad \mathcal{R}^k\begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} = \begin{bmatrix} a_1^{(1)} \\ a_2^{(1)} \\ \vdots \\ a_{n-k-1}^{(1)} \\ a_{n-k}^{(1)} \end{bmatrix} \tag{6}$$

20   where $\mathcal{R}$ denotes the RRO operator; and theoretically, we have $\mathcal{R}^k[a_1, a_2, \cdots, a_n] = [\hat{a}]$, if

21   $k \to \infty$. It clearly demonstrates that the RRO can continuously maintain the vital components

22   and reduce the dimensions of the original data. By continuously using the technique of low

23   rank estimation, Deep MF implements automatic estimation of dictionary size and number of

1    layers. Also, we provide a theoretical analysis of each operator and experimental validation of

2    Deep MF, in Sections 3 and 4 respectively, by comparing with three other deep linear models,

3    specifically, Deep SDL, Deep FICA and Deep NMF.

4

5    *2.3 Deep Sparse Dictionary Learning*

6

7         In the last decade, sparse dictionary learning (SDL), widely known as the algorithm Online

8    Dictionary Learning (ODL) (Mairal et al., 2010; Liu et al., 2010), has been successfully applied

9    to identify the concurrent BCNs of the human brain and the non-human primate brain from

10    fMRI datasets (Lv et al., 2015; Zhang et al., 2018). In this category, to satisfy the requirements

11    of hierarchical organization of BCNs, we propose a novel Deep SDL algorithm that is a

12    multilayer extension of conventional shallow SDL-based methods. Briefly, fMRI signals from

13    all voxels within the whole brain are extracted and are then organized as an extensive 2D

14    matrix, where the number of columns represents the total brain voxels and the number of rows

15    stands for the time points. For the first layer, the input 2D matrix is decomposed into the

16    product of an incomplete/over-complete dictionary basis matrix (each atom representing a

17    time series) and a feature matrix (representing this network's spatial volumetric distribution).

18    For each successive layer, the current features matrix is treated as an input matrix to be

19    continuously decomposed. A particularly important characteristic of this Deep SDL framework

20    is its ability to carry out over-complete decomposition for all layers; but, considering the finite

21    features, the deep layers only concentrate on the incomplete decomposition.

22

23         If considering all layers, the optimization function of Deep SDL is:

$$min_{D_i, \quad Y_i \in \mathbb{R}^{m \times n}} \left\| \prod_{i=1}^{M-1} D_i \, Y_M - SG \right\|_F^2 + [\mu_1 \ \mu_2 \cdots \mu_M] \otimes \sum_{i=1}^{M} \|Y_i\|_1 \qquad (7)$$

$$D_i Y_i \leftarrow Y_{i-1}$$

24    In Eq. (7), $\{D_i\}_{i=1}^{M}$ denotes the set of dictionaries; $i$ represents the number of current layer,

25    and $M$ is the total number of layers, $\{Y_i\}_{i=1}^{M}$ defines the set of hierarchical features, i.e., BCNs.

26    And $\{\mu_i\}_{i=1}^{M}$ represents all sparsity trade-offs for all layers, respectively.

27

Deep SDL, like ODL, utilizes GD as the optimizer to update all parameters. Simple GD is an efficient optimizer, but only guarantees the convergence of convex problems. In Section 3.1, we also provide the requirements of GD to be a contraction operator. Briefly, the convergence property of GD heavily depends on the step length. The following equation shows the iterative format of simple GD:

$$x_{k+1} \leftarrow x_k - \lambda f'_x \quad \lambda \in (0,1) \tag{8}$$

Compared with the ADMM optimizer used in Deep MF, the iterative function GD of Deep SDL is very simple. ADMM can utilize the current optima to update, and can be faster than GD, but ADMM requires that we can obtain the derivative of the optimization function.

*2.4 Deep Fast Independent Component Analysis*

ICA is a very popular and widely used data-driven computational technique, which was introduced to fMRI research over two decades ago (McKeown and Sejnowski 1998). In previous work, investigators have already reported that FICA using the Fixed-Point algorithm as an optimizer can be a very robust method (Hyvarinen, 1999). Inspired by Deep MF and FICA, the novel framework of Deep FICA aims to detect the hierarchically organized components. For simple shallow FICA, the original input matrix is decomposed as the weight matrix and independent component (IC) matrix. Applied to fMRI data, the ICs represent the BCNs. Similar to Deep MF and Deep SDL, in each layer of Deep FICA, the previous IC matrix is considered as the input signal matrix that will be decomposed using PCA and the Fixed-Point algorithm continuously (Figure 1). Deep FICA extracts only spatially independent features and can only solve the incomplete decomposition problem and not over-complete decomposition.

The optimization function of Deep FICA is:

$$max \sum_{k=1}^{M} \sum_{i=1}^{N} J_G(W_i^k) \tag{9}$$

In Eq. (9), $W_i^k$ represents the $i$th IC from the $k$th layer. The maximum $J_G(\cdot)$ indicates the independency of each potential IC.

1    Considering each IC, Deep FICA and FICA both utilize an efficacy Fixed-Point algorithm

2    to update the IC:

$$W_{it+1}^k \leftarrow \frac{W_{it}^k}{\sqrt{\|W_{it}^k C_{it}(W_{it}^k)^T\|}} \qquad (10)$$

3    Compared with Deep MF and Deep SDL, Deep FICA is relatively easy to implement, since

4    it does not include more complex algorithms, e.g., the RRO or the sparsity operator.

5

6    *2.5 Deep Non-negative Matrix Factorization*

7

8    Non-negative Matrix factorization (NMF) is a particularly useful family of techniques in

9    data analysis. Before the wide utilization of the Deep CNN, NMF was a crucial technique to

10    identify the features of a human face (Trigeorgis et al., 2016). In recent years, there has been

11    a significant amount of research on deep factorization methods that focus on particular

12    characteristics of both the data matrix and the hierarchical resulting factors. The application

13    area of the family of NMF algorithms has grown significantly during recent years. It has been

14    demonstrated that NMF can be a successful dimensionality reduction technique over a variety

15    of application areas including, but not limited to, environmetrics, microarray data analysis,

16    document clustering, face recognition and more. Moreover, due to its particular non-negative

17    constraints, NMF can also be directly utilized to analyze the fMRI data/signal (Lee & Seung,

18    1999). Deep NMF provides an opportunity to detect the potentially hierarchical structures of

19    BCNs.

20

21    Deep NMF focuses on the decomposition of the non-negative multivariate data matrix into

22    hierarchical factors $\{Z_i\}_{i=1}^M$ and $\{H_i\}_{i=1}^M$, such that (Trigeorgis et al., 2016):

$$
\begin{aligned}
Function_{Deep} &\overset{\text{def}}{=} \frac{1}{2}\|SG - Z_1 \cdot Z_2 \cdot \cdots \cdot Z_M H_M\|_F^2 \\
&= tr[SG^T SG - 2SG^T Z_1 \cdot Z_2 \cdot \cdots \cdot Z_M + SG_M^T \cdot Z_1 \cdot Z_2 \cdot \cdots \cdot Z_M \cdot Z_1^T \\
&\quad \cdot Z_2^T \cdot \cdots \cdot Z_M^T \cdot H_M]
\end{aligned}
\qquad (11)
$$

23    In Eq. (11), $SG$ represents the input fMRI signal, and $\{Z_i\}_{i=1}^M$ represents the weight

24    matrix; $\{H_i\}_{i=1}^M$ denotes the sets of non-negative components. $M$ denotes the total number

1    of layers. To calculate the optimal solutions of $\{Z_i\}_{i=1}^M$ and $\{H_i\}_{i=1}^M$ requires minimizing the

2    loss function $Function_{Deep}$ in Eq. (11). And $tr$ represents the trace of the matrix.

3

4        A key difference between Deep NMF versus Deep SDL and Deep MF is the updating

5    principle. Unlike Deep SDL and Deep MF, Deep NMF employs a fast policy to update these

6    two factors: $\{Z_i\}_{i=1}^M$ and $\{H_i\}_{i=1}^M$ (Trigeorgis et al., 2016). This principle is shown as follows:

$$H_{it+1}^k \leftarrow \frac{H_{it}^k}{f_{max}^H}$$
$$Z_{it+1}^k \leftarrow SG \cdot (H_{it+1}^k)^\dagger \tag{12}$$

7    where $H_{it}^k$ and $Z_{it}^k$ represent the non-negative components and weight matrix from the $k^{th}$

8    layer, iteration number *it*. And operator $(\cdot)^\dagger$ represents the pseudo-inverse of the input matrix

9    (Trigeorgis et al., 2016). The $f_{max}^H$ denotes the current maximum value of function *f*, related

10    to $H_{it}^k$.

11

12        Intuitively, these four deep linear models are each distinctive. Deep MF can be more

13    intelligent, and automatically determine all hyperparameters. Deep SDL can perform over-

14    complete decomposition for each layer. Deep FICA reveals the spatially independent and

15    hierarchical components, and has a faster convergence velocity. Finally, Deep NMF could

16    detect the non-negative components included in fMRI signals. We theoretically analyze the

17    relative performance of each deep linear model in the following section.

18

19

20           **3.      Results: Theoretical Analyses**

21

22        In this section, we employ mathematical theory, specifically real analysis, linear functional

23    analysis and abstract algebra, to explain why different deep linear models have the distinctive

24    characteristics that they do. In particular, we hope to explain:

25    (i) The advantages of linear deep models over shallow models and deep nonlinear models;

26    (ii) Why some deep linear models, e.g., Deep MF and Deep SDL, converge slowly while others,

27    e.g., Deep NMF and Deep FICA, converge quickly;

1  (iii) Why some deep linear models, e.g., Deep MF and Deep SDL, can better estimate

2  connectivity strength, while others, e.g., Deep NMF and Deep FICA, can better estimate the

3  spatial extent of connectivity networks.

4

5  *3. 1 Fundamental Interpretation of Each Linear Model*

6

7  All theoretical analyses are based on the vital assumption that all

8  mappings/operators/algorithms must be applied on a finite dimensional space. Please consult

9  Appendix A for the mathematical details (Assumption 1.1, Lemma 1.1, Theorem 1.1). If

10  considering any algorithm and/or process as an operator, Assumption 1.1 and Lemma 1.1

11  demonstrate that the norm of the operator should be equivalent, in order to dramatically

12  simplify our discussion.

13

14  According to Theorem 1.1, if considering the shallow linear, deep linear and deep

15  nonlinear models as approximations of the original function $f(x)$, then, obviously, deeper

16  models can employ more items such as $\{P_n(x)\}_{n=1}^{N}$ rather than just $P(x)$. Thus, the deeper

17  models can more accurately approximate the original function than a shallow model.

18  Meanwhile, nonlinear models can have: $\left\| \lim_{N\to\infty} \{P_n(x)\}_{n=1}^{N} - f(x) \right\| = 0$. But to optimize the

19  infinite items, it will be very time-consuming or even impossible to solve a non-polynomial (NP)

20  complexity problem. Hence, the theorem also answers why nonlinear models require a

21  sampling technique to reduce the complexity, e.g., Gibbs sampling for DBN (Hinton and

22  Salakhutdinov, 2006; Hinton et al., 2012).

23

24  According to the discussion in the last section, we can abstractly describe each deep

25  linear model using the combination of several operators. All operators involved in this study

26  are given in Table 1.

27

28

29

15

1

**Table 1**. All definitions of operators and their norms.

| Operator | Definition | Operator/Norm | Definition |
|---|---|---|---|
| $\mathfrak{U}$ | Deep MF | $\mathcal{U}$ | Update Operator of Deep NMF |
| $\mathfrak{N}$ | Deep NMF | $\mathcal{F}$ | Operator of Fixed-Point Algorithm |
| $\mathfrak{L}$ | Deep SDL | $\mathcal{C}$ | Consistent Operator |
| $\mathfrak{T}$ | Deep FICA | $\varphi$ | Norm of ADMM |
| $\mathcal{M}$ | Initialization Operator | $\rho$ | Norm of GD |
| $\mathcal{S}$ | Sparsity Operator | $\mu$ | Norm of Fixed-Point Algorithm |
| $\mathcal{P}$ | Principal Component Analysis (PCA) | $\gamma$ | Norm of Normalization in Deep NMF |
| $\mathcal{A}$ | ADMM | $\delta$ | Norm of Updating Deep NMF |
| $\mathcal{G}$ | GD | $C$ | Norm of Input fMRI Matrix |
| $\mathcal{N}$ | Normalization | $SG$ | Input fMRI Signal Matrix |
| $\mathcal{R}$ | Rank Reduction Operator | $\mathfrak{C}$ | Set of Consistent Operators |

2

3    In Table 2, we provide the definitions of sets involved in the following sections:

4    **Table 2**. All definitions of space and set

| Space/Set | Definition |
|---|---|
| $\mathbb{N}$ | Field of Natural Numbers |
| $\mathbb{R}$ | Field of Real Numbers |
| $\mathbb{K}$ | Field of Rational Numbers |

5

6    *3.2 Intensity Similarity*

7

8    As discussed in Sections 2.2 to 2.5, we have the following definitions:

9

10    **Definition 2.1** If we denote Deep MF as an operator $\mathfrak{U}$, based on the description of Deep MF,

11    considering the iteration $k$, we can denote $\mathfrak{U} \stackrel{\text{def}}{=} M \cdot \mathcal{A}^k \cdot \mathcal{S}^k \cdot \mathcal{R}^k$.

12    **Definition 2.2** If we denote Deep SDL as an operator $\mathfrak{L}$, based on the description of Deep

1     SDL, considering the iteration $k$, we can denote $\mathfrak{L} \overset{\text{def}}{=} M \cdot \mathcal{G}^k \cdot \mathcal{S}^k$.

2     **Definition 2.3** If we denote Deep FICA as an operator $\mathfrak{T}$, based on the description of Deep

3     FICA, considering the iteration $k$, we can denote $\mathfrak{T} \overset{\text{def}}{=} \mathcal{P} \cdot \mathcal{F}^k$.

4     **Definition 2.4** If we denote Deep NMF as an operator $\mathfrak{N}$, based on the description of Deep

5     NMF, considering the iteration $k$, we can denote $\mathfrak{N} \overset{\text{def}}{=} M \cdot \mathcal{U}^k \cdot \mathcal{N}$.

6

7     According to Definitions 2.1 to 2.4, as well as Corollaries 1.2 to 1.3 and Theorems 2.1 to

8     2.11 as proved in Appendix B, using the inequality of norm, considering the iteration $k$, and let

9     *SG* be the input matrix; for any operator applied on *SG*, we can derive the features as: $F_1$, $F_2$,

10     $F_3$, $F_4$; then we have:

$$\text{Deep MF:} \quad \|F_1\| = \left\|\mathfrak{A}^k \cdot SG\right\| \leq \|M\| \cdot \left\|\mathcal{A}^k\right\| \cdot \left\|\mathcal{S}^k\right\| \cdot \left\|\mathcal{R}^k\right\| \cdot \|SG\| \qquad (13\text{-}1)$$

$$\text{Deep SDL:} \quad \|F_2\| = \left\|\mathfrak{L}^k \cdot SG\right\| \leq \|M\| \cdot \left\|\mathcal{G}^k\right\| \cdot \left\|\mathcal{S}^k\right\| \cdot \|SG\| \qquad (13\text{-}2)$$

$$\text{Deep FICA:} \quad \|F_3\| = \left\|\mathfrak{T}^k \cdot SG\right\| \leq \left\|\mathcal{P}^k\right\| \cdot \left\|\mathcal{F}^k\right\| \cdot \|SG\| \qquad (13\text{-}3)$$

$$\text{Deep NMF:} \quad \|F_4\| = \left\|\mathfrak{N}^k \cdot SG\right\| \leq \|M\| \cdot \left\|\mathcal{U}^k\right\| \cdot \|\mathcal{N}\| \cdot \|SG\| \qquad (13\text{-}4)$$

11

12     According to Lemma 1.1 and Theorems 2.1 to 2.6, operators $\mathcal{A}^k$, $\mathcal{U}^k$, $\mathcal{N}$, $\mathcal{G}^k$, and $\mathcal{F}^k$

13     can be treated as contraction operators, which indicates that the norm of each operator should

14     be larger than zero and smaller than one. Other operators are constant values, according to

15     Theorems 2.1 to 2.11.

16

17     If we denote the norm of contraction operators as:

$$\left\|\mathcal{A}^k\right\| = \varphi < 1 \qquad (14\text{-}1)$$

$$\left\|\mathcal{G}^k\right\| = \rho < 1 \qquad (14\text{-}2)$$

$$\left\|\mathcal{F}^k\right\| = \mu < 1 \qquad (14\text{-}3)$$

$$\left\|\mathcal{U}^k\right\| = \delta < 1; \quad \|\mathcal{N}\| = \gamma < 1 \qquad (14\text{-}4)$$

18     Despite the fact that the norms of operators $\mathcal{A}, \mathcal{U}, \mathcal{G}, \mathcal{F}$ are not equivalent, according to

19     Theorems 3.1 and 3.2 (see Appendix C), we consider an extreme condition $k \to \infty$, and then

20     we have: $\varphi = \delta = \rho = \mu$ that indicates convergence to the global optimum. Then we can

21     rewrite all equations 13-1 to 13-4 as:

$$Deep\ MF: \quad \|F_1\| \leq \quad \varphi \cdot \|M\| \cdot \|\mathcal{S}^k\| \cdot \|\mathcal{R}^k\| \cdot C \tag{15-1}$$

$$Deep\ SDL: \quad \|F_2\| \leq \rho \cdot \|M\| \cdot \|\mathcal{S}^k\| \cdot C \tag{15-2}$$

$$Deep\ FICA: \quad \|F_3\| \leq \quad \mu \cdot \|\mathcal{P}^k\| \cdot C \tag{15-3}$$

$$Deep\ NMF: \quad \|F_4\| \leq \delta \cdot \gamma \cdot \|M\| \cdot C \tag{15-4}$$

Obviously, based on Eqs. (15-1) to (15-4), we have the conclusion:

$$\|F_4\| \leq \|F_3\| \leq \|F_2\| \leq \|F_1\| \tag{16}$$

Since all features $\{F_i\}_{i=1}^4$ have the same dimensions, this inequality Eq. (16) can clearly explain why the intensity of features, i.e., the connectivity strength of voxels in the networks, varies based on the different models. In particular, $F_4$ (the features derived from Deep NMF) should have the smallest intensity and $F_1$ (the features obtained via Deep MF) should have the largest intensity. Meanwhile, Eqs. (15-1 to 15-4) also reveal the convergence velocity of each model. Since Deep MF contains the most operators with the complex optimization function ADMM, it should be slowest. Because Deep SDL uses a sparsity operator as well as GD, which is relatively slow, it is comparable in speed to Deep MF, even given a perfect step-length. Theoretically, Deep FICA and Deep NMF should have faster convergence.

*3.3 Spatial Similarity*

Spatial matching is another important way to measure the similarity between identified components and templates. To examine this property, we use Assumptions 3.1 to 3.2 and Lemma 3.1 to prove Theorem 3.1 in Appendix C:

**Theorem 3.1** If we denote the following sets:

$$Deep\ MF: D = \{\mathfrak{A}^k N, N \in \bigcup_{i=1}^{M} voxel_i,\ voxel_i \notin T\}$$

$$Deep\ SDL: L = \{\mathfrak{L}^k N, N \in \bigcup_{i=1}^{M} voxel_i,\ voxel_i \notin T\}$$

$$Deep\ FICA: I = \{\mathfrak{T}^k N, N \in \bigcup_{i=1}^{M} voxel_i,\ voxel_i \notin T\}$$

$$Deep\ NMF: \Theta = \{\mathfrak{R}^k N, N \in \bigcup_{i=1}^{M} voxel_i,\ voxel_i \notin T\}$$

1    And considering the iteration $k$, and $k > K$, it implies:

$$0 < \frac{|\mathfrak{A}^k V|}{|V \cup D|} \leq \frac{|\mathfrak{L}^k V|}{|V \cup L|} \leq \frac{|\mathfrak{T}^k V|}{|V \cup I|} \leq \frac{|\mathfrak{N}^k V|}{|V \cup \Theta|} \qquad (17)$$

2    where $|\cdot|$ denotes the number of positive elements. $\Theta$ represents the set that only contains

3    the element $0$. $T$ represents the functional regions of brain.

4

5    Since the convergence of deep models is a vital issue when solving real world problems

6    (Topol, 2019), Eq. (16) and Theorem 3.1 can explain the convergence of all deep linear models,

7    considering enough iterations. Clearly, if we examine the spatial similarity between two BCNs,

8    according to Theorem 3.1, we can conclude: with the same number of iterations, Deep NMF

9    has the best performance on spatial matching, Deep FICA has the next best performance, and

10   Deep SDL and Deep MF have the least. That is, the norm of the operator of Deep NMF is very

11   small, and is iteratively applied on functional regions and background noise, which causes the

12   intensity of functional areas to decrease very rapidly. However, since the intensity of

13   background is very small, the noise can be reduced to near zero much faster than Deep MF

14   and Deep SDL. The performance of Deep FICA on spatial matching should be comparable to

15   Deep NMF; since a normalization operator is involved in Deep NMF, the intensity of

16   components identified by Deep NMF should be smaller than Deep FICA. Theorem 3.2

17   included in Appendix C also explains that all deep linear models finally converge given enough

18   iterations.

19

20   To test these theoretical analyses, in the next section, a simulated experimental

21   reconstruction will be introduced as the ground truth templates for the first layer BCNs. These

22   templates will be employed to construct the simulated fMRI signal and all deep linear models

23   will be applied on the simulated data and their $1^{st}$ layer results will be compared to the

24   templates. By examining the intensity similarity and spatial similarity to the ground truth

25   templates, the correctness of the theoretical conclusions can be investigated.

26

27

# 4. Results: Experimental Validation

*4.1 Simulated fMRI Data*

In this work, we employ an *in silico* fMRI simulation method proposed previously (Zhang et al., 2018, 2019), using templates of BCNs (Smith et al., 2009) to test these proposed deep linear models. Specifically, we selected 12 BCNs (Table 3) that were originally derived using conventional shallow ICA (Smith et al., 2009).

**Table 3**. All abbreviations of BCNs in simulation

| Name/Number | Abbreviation | Name | Abbreviation |
|---|---|---|---|
| Primary Visual Network/1 | VIS-1 | Auditory Network/7 | AUD |
| Perception Visual Shape Network/2 | VIS-2 | Executive Control Network/8 | ECN |
| Perception Visual Motion Network/3 | VIS-3 | Left Frontoparietal Network/9 | FP-L |
| Default Mode Network/4 | DMN | Right Frontoparietal Network/10 | FP-R |
| Brainstem & Cerebellum Network/5 | B/C | Dorsal Attention Network/11 | DAN |
| Sensorimotor Network/6 | SM | Salience Network/12 | SN |

These template BCNs derived from resting-state fMRI have been released publicly and are considered to be functional brain areas covering a large part of cerebral cortex (Smith et al., 2009). Since all deep linear models should be evaluated equally using a known ground truth, we employ a simulation of resting-state fMRI signals. Following the fMRI simulation pipeline for spatially independent networks named Experiment 1 in the previous study (Zhang et al., 2019), the 12 templates (Smith et al., 2009) are collected as components/spatial features; and we adopt 12 time series, including 200 time points, derived from a previous study

1    (Lv et al., 2015). The final simulation data is a matrix obtained as the product of time series

2    and components.

3

4      The detailed parameters of each template are: 91×109 matrix, 91 slices, 2.0 mm isotropic

5    voxels. The number of mask voxels is 262,309 and the number of time points is 200. All

6    templates are registered to standard MNI space at 2.0 mm. This pipeline contains the steps

7    of spatial artifact cleanup, distortions removal and cortical surfaces generation. After that,

8    different subjects are aligned to the standard MNI space (Lv et al., 2015; Zhang et al., 2019).

9

10      Table 4 provides the main hyperparameter settings of the four proposed deep linear

11    models, including: the number of components of the $1^{st}$ and $2^{nd}$ layers, the number of iterations

12    and the step length of gradient descent, where applicable. Since the 1st layer may include

13    noise components, we choose a larger number of components than the expected number of

14    features, which in this case is at least a dozen ground truth template BCNs. For the 2nd layer,

15    which should have fewer high-level features, the number of components should be less than

16    in the $1^{st}$ layer. Lv et al. (2015) introduce an experimental method to search for the best number

17    of components, but, in fact, heuristically tuning the hyperparameters of deep models is very

18    difficult. Since Deep MF is capable of estimating all these hyperparameters automatically, only

19    the maximum number of iterations is given. For the other three methods, the hyperparameter

20    values were chosen heuristically based on best matching to the ground truth templates for the

21    $1^{st}$ layer and for perceived quality of the derived networks for higher layers.

22

23    **Table 4**. Important Hyperparameter Settings of Four Deep Linear Models

|  | Deep MF | Deep SDL | Deep FICA | Deep NMF |
|---|---|---|---|---|
| Number of Components of $1^{st}$ layer | N/A | 15 | 20 | 30 |
| Number of Components of $2^{nd}$ layer | N/A | 13 | 10 | 15 |
| Number of Iterations | 100 | 100 | 100 | 100 |
| Step Length | N/A | 0.01 | N/A | N/A |

24

1    *4.2 Investigating the First Layer Reconstructions of Each Deep Linear Model via Intensity,*

2    *Spatial and Hausdorff Distances*

3

4    We can quantitatively compare the identified components, i.e., BCNs, with the original

5    ground truth, i.e., templates, in three distinct ways. First, the similarity can be calculated

6    spatially, largely independent of the intensity of each voxel of the identified components. The

7    definition of spatial similarity is:

$$Similarity Spatial = \frac{|Component \cap Template|}{|Component \cup Template|} \quad\quad (18)$$

8    where $|\cdot|$ represents binarization, which represent the voxels above a given intensity

9    threshold. The spatial similarity is measuring the ratio of intersection and union of identified

10    component and template.

11    In contradistinction, only considering the intensity of each voxel of the derived

12    components, it is useful to calculate the distance between the intensities of components and

13    templates. The definition of intensity similarity is:

$$Similarity_{Intensity} = \left( \sum_{i=1}^{N} \frac{|x_i - y_i|}{|x_i| + |y_i|} \right)^{-1} \quad\quad (19)$$

14    where $|\cdot|$ represents the absolute value. Given a threshold, the intensity similarity is

15    calculated via summed absolute value of intensity of component (denoted as $x_i$) and template

16    (denoted as $y_i$) divided by the absolute value of their difference. $N$ denotes the total number

17    of voxels. Obviously, if all intensity values of identified component and template are equal, the

18    intensity similarity approaches infinity.

19

20    Finally, to jointly consider both spatial and intensity matching, we use the Hausdorff
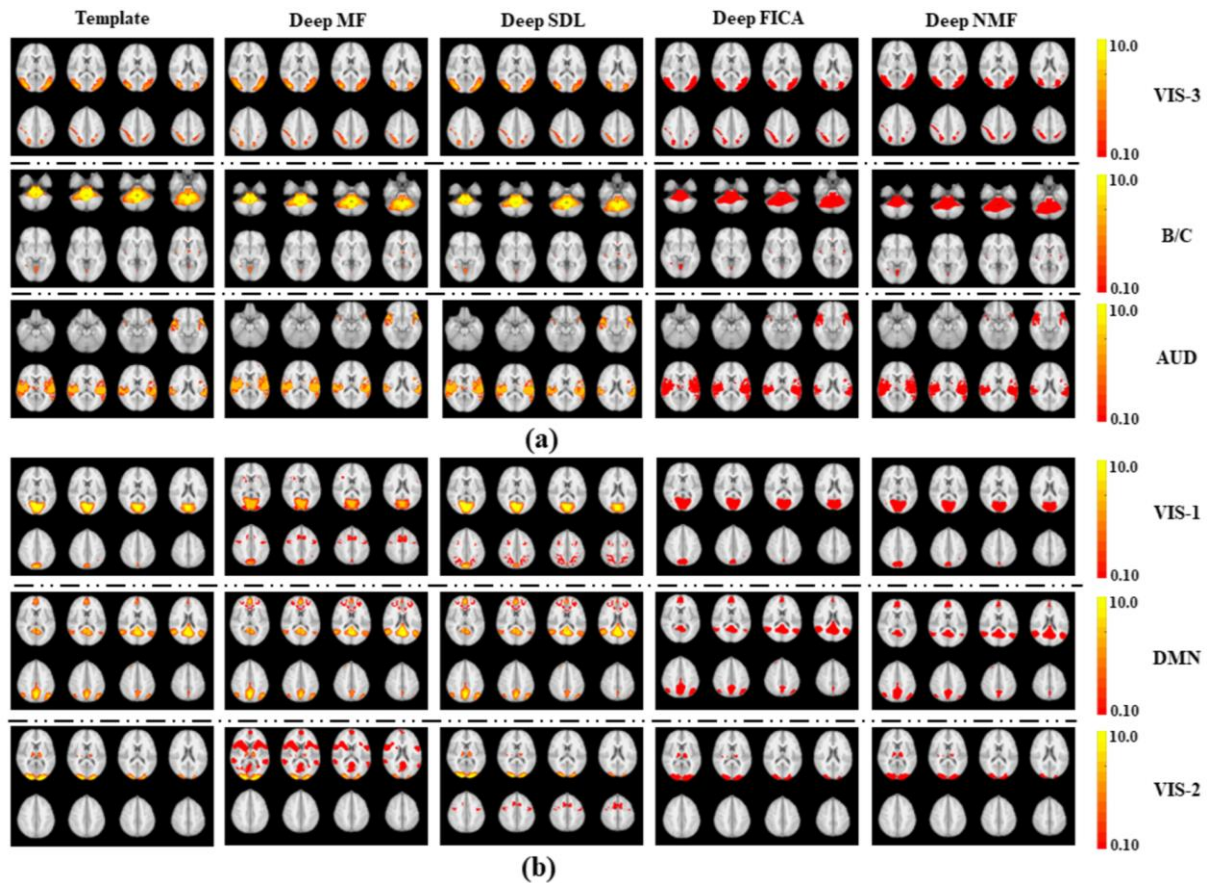
21    Distance (HD):

$$X = \sum_{i=1}^{N} 2 \times min(x_i, y_i), if\ x_i,\ y_i \in C \cap T$$

$$Y = \sum_{i=1}^{M} x_i + y_i, if\ x_i,\ y_i \in C \cup T \quad\quad (20)$$

$$HD = \frac{X}{Y}$$

22    Briefly, $X$ represents two times the minimum intensity value of intersection of component

1 and template; and $Y$ represents the summed intensity value of union of component and

2 template. $C$ and $T$ represent the sets of components and templates, respectively. Therefore,

3 HD includes the influences of intensity similarity and spatial overlap simultaneously.

4



5

6 **Figure 2.** Comparison of six 1st layer networks from all four deep linear models with the ground truth
7 templates from simulated fMRI data (see Table 3 for network abbreviations). The first column presents
8 eight representative slices from each of six representative template networks. The second to fifth
9 columns show the corresponding slices from the networks identified via Deep MF, Deep SDL, Deep
10 FICA and Deep NMF, respectively. **(a)** The AUD, B/C and VIS-3 networks illustrate better intensity
11 matching to the templates by Deep MF and Deep SDL than by Deep FICA or Deep NMF (see color bar
12 of intensities measuring connectivity strength on the right). **(b)** The VIS-1, VIS-2 and DMN networks
13 also show this same disparity among the deep linear models for intensity matching, but also show better
14 spatial similarity to the templates for Deep FICA and Deep NMF compared to Deep MF or Deep SDL.

15

16 The results show that Deep NMF and Deep FICA produce smaller network intensities than

17 the templates, whereas Deep MF and Deep SDL yield larger intensities that better match the

18 templates (Figure 2a). In contradistinction, there are generally more noisy areas detected from

19 Deep MF and Deep SDL due to the larger norms of their iterative operators, compared to Deep

NMF and Deep FICA (Figure 2b). Hence, Deep NMF and Deep FICA have better spatial similarity to the templates than the other two methods. This illustrates the trade-off between intensity matching and spatial matching. To view reconstructions for all 12 examined BCNs and for further details, please see Figure S1 included in the Supplemental Materials.

As defined by Eq. (18) to Eq. (20) in Section 4.2, the quantitative comparisons among the four deep linear models for intensity similarity, spatial similarity and the Hausdorff distance are provided by Figure 3. These quantitative results clearly demonstrate that Deep MF and Deep SDL provide the best intensity matching (Figure 3a), since their convergence velocity is relatively slow. Therefore, Deep MF and Deep SDL can reconstruct the most accurate connectivity strengths of each component from input fMRI signals, consistent with theory (Section 3). Considering spatial similarity, due to the fastest convergence velocity and non-negative normalization of Deep NMF, the intensity is reduced rapidly across iterations. Since the noise has smaller intensity than the signal, it is reduced much faster, which helps account for Deep NMF yielding the best spatial similarity results for most networks (Figure 3b). This result is also predicted by theory in Section 3. A rigorous proof is presented in the Appendices. It should be noted that Deep FICA has an inherent advantage over the other models for spatial matching since it is most similar to the shallow ICA analysis used to generate the ground truth templates for the BCNs.
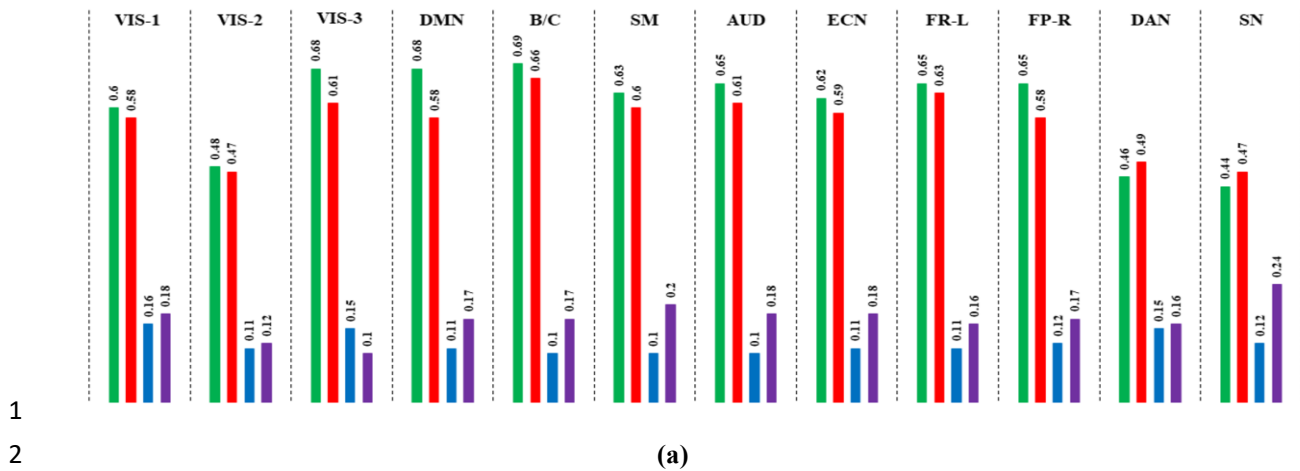
All proposed deep linear models can be evaluated by HD to consider both intensity and spatial similarity (Figure 3c). Deep MF generated the best performance for all 12 BCNs with Deep SDL running close behind. Hence, the additional RRO in Deep MF does yield advantages over the other three deep linear models. Similarly, the sparsity operator of Deep SDL and Deep MF help them outperform Deep FICA and Deep NMF, which both lack that capability.
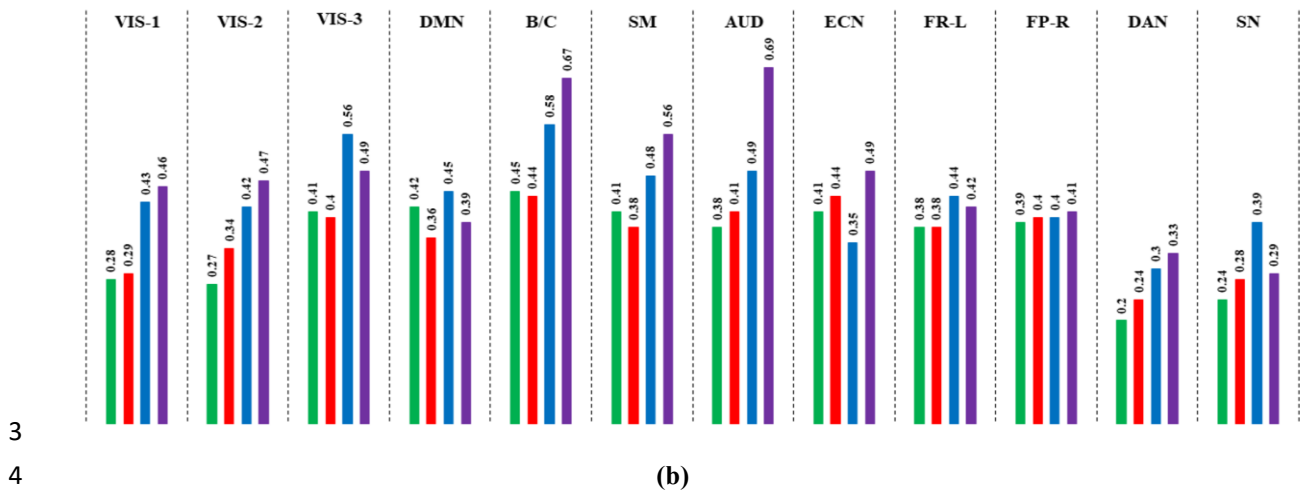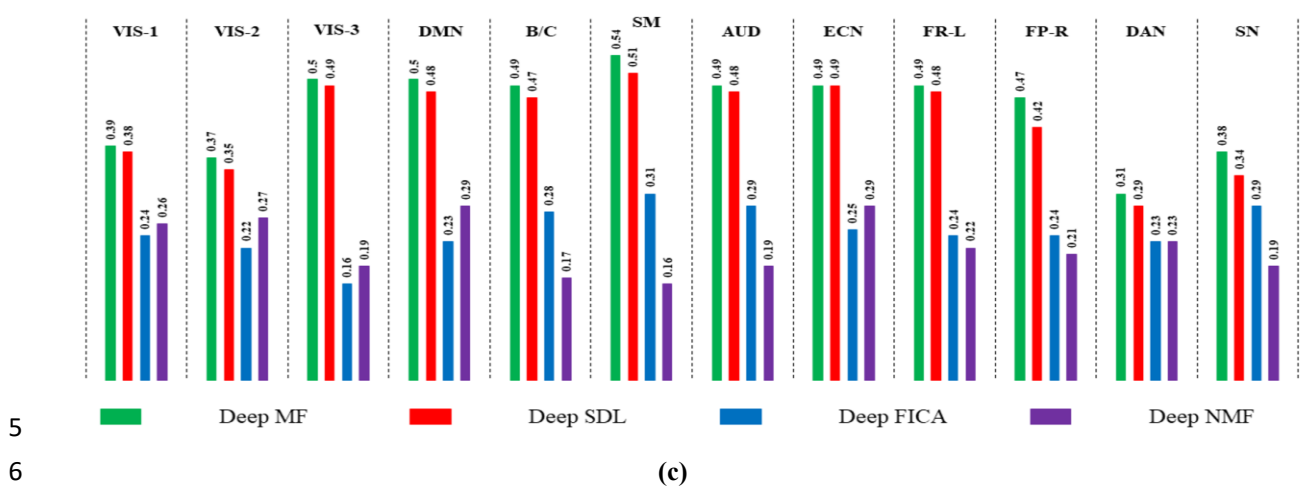
**Figure 3.** Comparisons of the twelve 1st layer networks of the four deep linear models for **(a)** intensity similarity to the ground truth templates; **(b)** spatial similarity to the ground truth templates; and **(c)** the Hausdorff Distance to the ground truth templates that jointly considers intensity and spatial similarity.

1   *4.3 The 2$^{nd}$ Layer Networks of Each Deep Linear Model*

2

3   Compared with the shallow 1$^{st}$ layer features, it is difficult to successfully investigate the

4   features of deeper layers because there is no widely accepted ground truth for those more

5   complex higher-level networks. The 2$^{nd}$ layer features can be comprehended as the

6   recombination of 1$^{st}$ layer features. Another challenge for testing deeper networks is that the

7   deep linear models differ with regard to how many layers can be reconstructed from a given

8   dataset. For example, Deep FICA can only decompose the simulated fMRI into two layers, but

9   Deep MF can decompose the simulated signal into four layers. Given these constraints, as

10  well as the altered connectivity strengths in the 2$^{nd}$ layer relative to the 1$^{st}$ layer, we limit the

11  analysis of deeper networks to examining the spatial similarity between 2$^{nd}$ layer networks of

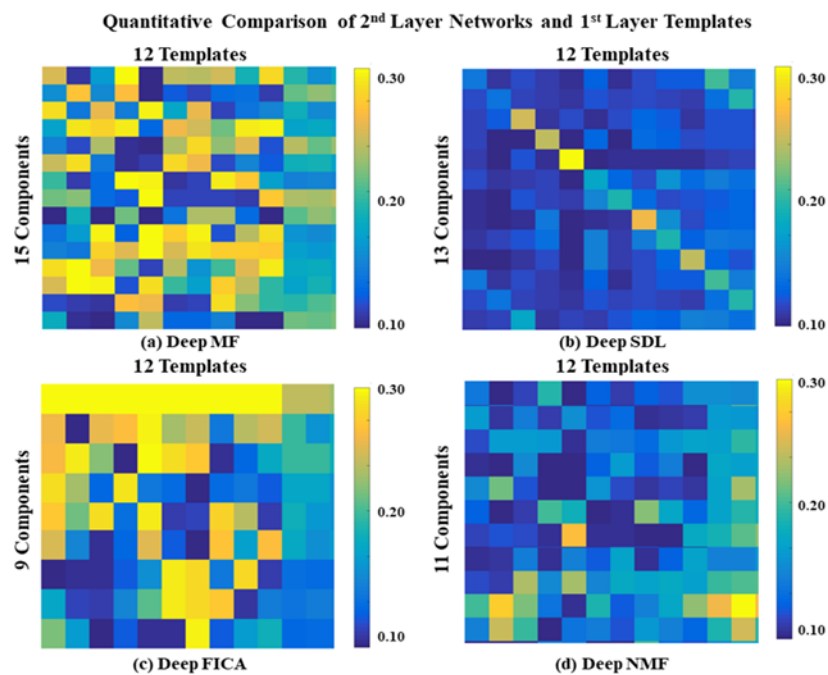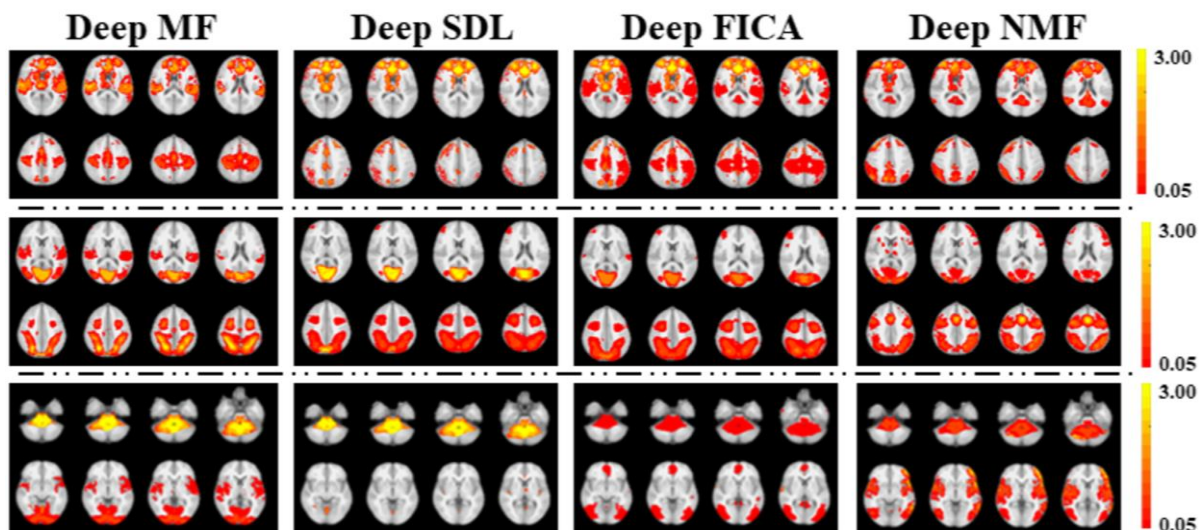12  each of the four deep linear models with the shallow ground truth templates (Figure 4).

13



14

15  **Figure 4.** Comparisons of BCNs from 2$^{nd}$ layer of Deep MF, Deep SDL, Deep FICA and Deep NMF.

16  Each element represents the spatial similarity of the identified component and the ground truth

17  templates; (a), (b), (c) and (d) are Deep MF, Deep SDL, Deep FICA and Deep NMF, respectively. The

18  rows represent the identified 2$^{nd}$ layer BCNs and the columns represent the ground truth templates of

19  the simulated experiment.

20

21

22

23

1    Based on these preliminary comparisons, it is clear that the four deep linear models can

2    produce different higher-level features. Most notable is that Deep SDL produces $2^{nd}$ layer

3    networks that are the most spatially similar to the shallow ground truth templates, as shown

4    by the larger main diagonal elements in its similarity matrix and the smaller off-diagonal

5    elements (Figure 4b). Hence, Deep SDL does relatively little recombination of the $1^{st}$ layer

6    features in its $2^{nd}$ layer. In contradistinction, the first component of the Deep FICA $2^{nd}$ layer

7    (top row of its similarity matrix in Figure 4c) is very strongly correlated with 10 of the 12 ground

8    truth templates and therefore appears to be a spatially "global" network. The $2^{nd}$ layer features

9    of Deep NMF have overall the least spatial similarity with the ground truth templates (Figure

10   4d) whereas Deep MF produces the greatest variation in the correlations between its $2^{nd}$ layer

11   features and the ground truth templates (Figure 4a).

12



13

14   **Figure 5.** Comparisons of BCNs derived from the $2^{nd}$ layer of Deep MF, Deep SDL, Deep FICA and
15   Deep NMF. Each column includes three representative $2^{nd}$ layer networks from a deep linear model,
16   matched across models in each row.

17

18   Three representative $2^{nd}$ layer BCNs matched for each deep linear model are presented

19   in Figure 5. The full set of non-noise $2^{nd}$ layer networks are given in Figure S2 of the

20   Supplemental Materials. The top row of Figure 5 shows that the nodes of the ECN, including

21   anterior cingulate cortex and medial prefrontal cortex, are represented in that $2^{nd}$ layer network

22   for all four models. For both Deep MF and Deep FICA, the ECN is combined with nodes of the

23   SN, including the insulae, pre-supplementary motor areas (pre-SMA), and premotor areas.

1    For Deep NMF and Deep SDL, however, the ECN is joined instead with nodes of the DMN,

2    including the precuneus, posterior cingulate cortex and the superior parietal lobules. This

3    higher-level connectivity between ECN and DMN is weaker for Deep SDL than Deep NMF,

4    whereas connectivity within ECN is stronger for Deep SDL than Deep NMF, in keeping with

5    the observations from Figure 4 that Deep SDL preserves $1^{st}$ layer networks the most of all four

6    algorithms, whereas Deep NMF preserves $1^{st}$ layer networks the least (Figures 4 & S2). It can

7    be observed that the same $2^{nd}$ layer network of Deep FICA also contains parts of the DMN,

8    most notably the posterior cingulate cortex, although to a lesser extent than Deep NMF.

9    Therefore, Deep FICA recombines nodes of three different $1^{st}$ layer spatially independent

10    components (DMN, ECN & SN) into a single $2^{nd}$ layer independent component. Figure S3 of

11    the Supplementary Materials provides a spatial similarity matrix for the non-noise $2^{nd}$ layer

12    networks for Deep SDL, Deep FICA and Deep NMF with reference to Deep MF.
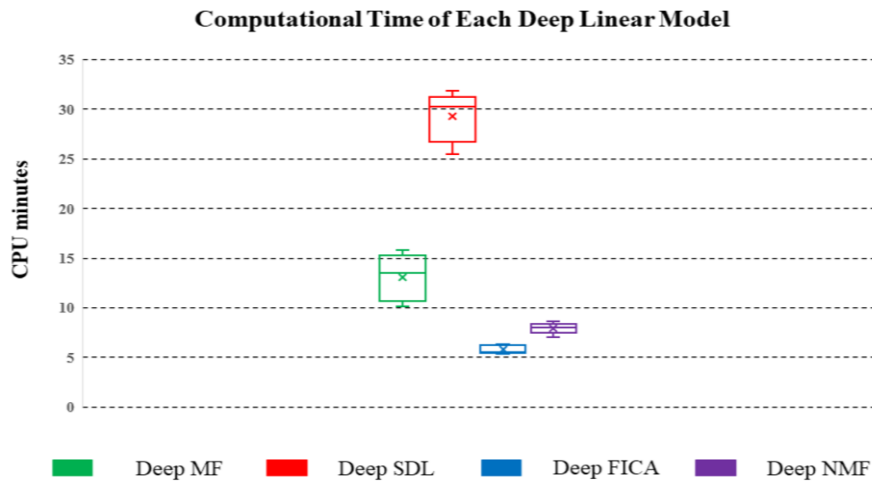
13

14    The middle row of Figure 5 shows that all four models produced a $2^{nd}$ layer network

15    consisting of VIS-1 and DAN. However, Deep NMF additionally included parts of VIS-2, VIS-

16    3 and pre-SMA whereas Deep MF additionally included VIS-3 and posterior perisylvian regions.

17    The bottom row of Figure 5 illustrates links of the brainstem and cerebellum with visual cortex.

18    However, Deep MF finds correlations of B/C with VIS-1 & VIS-2 whereas Deep FICA and Deep

19    NMF finds correlations with VIS-3 instead. Both Deep MF and Deep NMF include perisylvian

20    regions in this $2^{nd}$ layer network as well. Given the relatively slow convergence of Deep SDL

21    via gradient descent, these features might be expected in its $3^{rd}$ layer instead.

22

23    We also compare computation time for the four deep linear models presented in this work

24    on our computing cluster (Figure 6), which demonstrates that Deep FICA is the fastest and

25    Deep SDL is by far the slowest. Deep MF provides the best trade-off between speed and

26    performance as judged by reconstruction accuracy for $1^{st}$ layer networks (Figure 3).

**Figure 6.** Comparisons of computation time of 10 independent runs, using the same number of iterations and the same simulated fMRI dataset for Deep FICA (blue), Deep SDL (red), Deep MF (green) and Deep NMF (purple). The box plots give the mean and standard deviation of the CPU time in minutes for the 10 runs.


# 5. Discussion

We have introduced novel deep linear models that integrate multiple operators to extract hierarchical spatial features in fMRI data. These models bridge the gap between traditional shallow linear models (Andersen et al., 1999; Beckmann et al., 2005; Calhoun et al., 2001; Hyvarinen, 1999; Lee et al., 2011; Lee et al., 2016; Mairal et al., 2010; McKeown & Sejnowski, 1999) and newer deep nonlinear models (Hu et al., 2018; Huang et al., 2018; Dong et al., 2020; Zhang et al., 2020). The primary advantages of the proposed algorithms over more complex deep nonlinear models are to quickly and easily map the hierarchical organization of BCNs without requiring large amounts of fMRI data or HPC clusters with GPUs or TPUs. The behavior of deep linear models is also more explainable than are, for example, CNNs and DBNs, as we show through theoretical predictions of their relative performance (Section 3) that are validated via simulations (Section 4). Furthermore, convergence to the global optimum can be guaranteed for deep linear models with convex optimization functions, unlike deep nonlinear models where such convergence is rarely achieved in practice. This is important given the recent realization that real-world imaging applications often suffer from underspecification, resulting in wildly unpredictable performance from any particular deep nonlinear network due to convergence to different local optima from different random initial conditions despite identical training data and hyperparameters (D'Amour et al., 2020).

1    Deep MF employs ADMM, which is a distributed optimization algorithm particularly well
2    suited to compositional analysis of hierarchical modular systems, and also utilizes RRO for
3    data-driven determination of all hyperparameters, which can be considered an intelligent
4    factorization method. This is a major advantage over many conventional shallow data-driven
5    fMRI connectivity reconstruction methods and the other three deep linear models presented
6    here, as well as more complex deep nonlinear models, all of which must be manually tuned
7    for hyperparameter settings. Deep SDL can explore more potential components than other
8    models, even more than the number of original time points, via over-complete decomposition.
9    Deep NMF converges very rapidly and recognizes the non-negative constraints of BCNs in
10   fMRI. Finally, Deep FICA efficiently maps hierarchically spatially independent BCNs and is
11   even easier to implement than the other peer deep linear models, especially given the wide
12   usage of shallow ICA models for unsupervised fMRI mapping.

13

14   In this research, we also introduce an innovative framework for studying the relative
15   performance of deep linear models, both theoretically by comparing their mathematical
16   structure as well as *in silico* via fMRI simulations. Evaluating the $1^{st}$ layer reconstructions of
17   these deep linear models using simulated fMRI data from widely accepted ground truth BCNs,
18   we find that Deep MF and Deep SDL are clearly superior for computing connectivity strength
19   whereas Deep NMF and Deep FICA are modestly better for mapping spatial extent. These
20   results were predicted from the unique mix of mathematical operators used in each of the four
21   methods (Section 3). Overall, Deep MF provided the most robust combination of intensity
22   matching, spatial matching and computational efficiency of all four techniques. This can be
23   attributed to its joint use of sparsity and rank reduction operators in conjunction with the
24   distributed ADMM optimization function. We also discovered that deeper features such as the
25   $2^{nd}$ layer BCNs are recombinations of the $1^{st}$ layer networks and that these can vary among
26   the four deep linear models. For example, Deep SDL produces the least recombination of the
27   $1^{st}$ layer networks in its $2^{nd}$ layer. This can be attributed to its relatively slow convergence
28   velocity using gradient descent optimization; therefore, more low-level network recombination
29   is seen in its $3^{rd}$ layer instead. Another important factor that differs among the deep linear
30   models is the number of spatial features that can be accommodated at each level of the

1    hierarchy and the maximum number of layers for any given dataset. For example, Deep FICA

2    supports the fewest number of meaningful components at the 2$^{nd}$ layer (nine) and only two

3    layers total; therefore, its 2$^{nd}$ layer networks would be the least sparse. This can be seen in

4    the top row of Figure 5 in which Deep FICA combines nodes of DMN, ECN and SN into a

5    single 2$^{nd}$ layer network, unlike the other models that only incorporate two of the three 1$^{st}$ layer

6    networks, but which can instead generate even deeper networks beyond the 2$^{nd}$ layer. Hence,

7    the choice of deep linear model matters for exploring higher-level BCNs. The mathematical

8    evaluation framework and the fMRI simulation procedure provided in this work should enable

9    further development of deep linear models that are optimized for different types of real-world

10    applications in biomedical imaging, with Deep MF as the current best algorithm for fMRI

11    hierarchical functional connectivity mapping.

12

13    One shortcoming of the current work is that the ground truth templates for testing the 1$^{st}$

14    layer networks were generated using conventional shallow ICA (Smith et al., 2009), which is

15    currently the most widely accepted technique for data-driven analysis of functional connectivity.

16    Aside from giving Deep FICA an inherent advantage, these spatially independent BCNs do

17    not adequately evaluate the ability to reconstruct overlapping networks that is a property of

18    methods such as shallow or Deep SDL. We also do not comprehensively investigate the

19    properties of the deeper layers of these four models, which is an extensive topic that is beyond

20    the scope of this paper, especially considering the absence of gold standards for these more

21    complex high-level networks as well as the wide variation among deep linear models in key

22    attributes such as convergence velocity and enforcement of sparsity.

23

24    In this initial exploratory work, many of the derived 2$^{nd}$ layer BCNs demonstrate

25    neurobiological face validity. For example, the SN is known to modulate the anticorrelated

26    connectivity of the DMN and the ECN (Menon & Toga, 2015); hence the linkage of their nodes

27    into a single higher-level network (Figure 5, top row). The functional coupling of vision

28    networks with the DAN shown in Figure 5 (middle row) is also well known, given the role that

29    the DAN plays in visual attention and eye movements (Vossel et al., 2014). Future

30    neuroscientific studies will be required to empirically validate the deep features of these

1    models using demographic, clinical, cognitive, behavioral and/or electrophysiological data.

2

3    Since these deep linear models do not require large training datasets nor specialized

4    computing infrastructure, they can be easily applied to clinical research with the potential to

5    generate novel functional connectivity biomarkers of neurodevelopmental, neurodegenerative,

6    and psychiatric disorders (Parkes et al., 2020), including for diagnosis, prognosis and

7    treatment monitoring. This is particularly significant given the recent observation that

8    neuropathology and psychopathology often affect low-level network connectivity differently

9    than high-level network connectivity. For example, many different psychiatric disorders have

10   been found to decrease lower-order sensory and somatomotor network connectivity in a

11   uniform manner across patients (Elliott et al., 2018; Kebets et al., 2019), while increasing

12   distinctiveness among patients in networks at higher levels of the hierarchy (Kauffman et al.,

13   2017; Parkes et al., 2020). In fMRI studies of mild traumatic brain injury (TBI), altered

14   functional connectivity has been found early after concussion both within individual BCNs,

15   such as the SN, DMN and ECN, as well as between different BCNs (Palacios et al., 2017).

16   Interactions of BCNs, such as that of the SN with the DMN, are thought to be especially

17   important for outcome after TBI and can be used to guide personalized treatment (Jilka et al.,

18   2014; Li et al., 2019). Disordered coupling of the SN with the DMN and ECN has also been

19   shown in mild cognitive impairment (Chand et al., 2017). Hence, prevalent neurological

20   disorders such as head trauma and neurodegenerative disease are thought to affect multiple

21   levels of the human brain's hierarchical organization. Such high-level interactions between

22   DMN, ECN and SN can be investigated with deeper layers of these hierarchical linear models

23   that integrate their spatially distinct gray matter nodes into a single larger-scale network, as

24   seen in Figure 5 (top row). These examples show how more principled data-driven

25   characterization of this hierarchy, particularly at its higher levels, holds great promise for

26   providing clinically actionable biomarkers of neurological and psychiatric diseases.

27

28   The benefits of deep linear models gain in importance as the spatial and temporal

29   resolution and sensitivity of fMRI continue to increase with improved MR imaging hardware

30   and pulse sequences, e.g., the advent of SLice Dithered Enhanced Resolution Simultaneous

1    MultiSlice (SLIDER-SMS) imaging (Vu et al., 2018) and MultiBand MultiEcho (MBME) imaging

2    (Boyacioğlu et al., 2015; Cohen et al., 2020). Higher fMRI sensitivity and spatial resolution will

3    enable mesoscale functional imaging that supports more 1$^{st}$ layer components of deep linear

4    models to uncover subnetworks of the BCN templates used in this work. This will also permit

5    the use of deeper models to extract more levels of the hierarchy of functional connectivity.

6    Whereas many widely used methods for performing time-varying fMRI analysis are heuristic

7    rather than data-driven, such as those with arbitrary time windows (Iraji et al., 2020), advances

8    in fMRI temporal resolution can be combined with deep linear models that perform joint

9    spatiotemporal decomposition for principled unsupervised dynamic functional connectivity

10    mapping that reveals ever more of the human brain's hierarchical organization.

11

## 6. Acknowledgements

15

# References

Andersen, A.H., Gash, D.M., Avison, M.J. (1999). Principal component analysis of the dynamic response measured by fMRI: a generalized linear systems framework. *Magnetic Resonance Imaging*, 17:795-815.

Bartels, A., Zeki, S. (2005). Brain dynamics during natural viewing conditions - a new guide for mapping connectivity in vivo. *Neuroimage*, 24:339-349.

Bassett, D. S., Bullmore, E., Verchinski, B. A., Mattay, V. S., Weinberger, D. R., & Meyer-Lindenberg, A. (2008). Hierarchical organization of human cortical networks in health and schizophrenia. *Journal of Neuroscience*, 28:9239-9248.

Beck, A., & Teboulle, M. (2009). A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM Journal on Imaging Sciences*, 2:183-202.

Beckmann, C.F., Smith, S.M. (2005). Tensorial extensions of independent component analysis for multisubject FMRI analysis. *Neuroimage*, 25:294-311.

Bengio, Y., Courville, A.C., Vincent, P. (2012). Unsupervised feature learning and deep learning: A review and new perspectives. CoRR, abs/1206.5538, 1.

Biswal, B., Yetkin, F.Z., Haughton, V.M., Hyde, J.S. (1995). Functional connectivity in the motor cortex of resting human brain using echo-planar MRI. *Magn Reson Med*, 34(4):537-41.

Biswal, B.B., Maarten, M., Xi-Nian, Z., Suril, G., Clare, K., Smith, S.M., Beckmann, C.F., Adelstein, J.S., Buckner, R.L., Stan, C. (2010). Toward discovery science of human brain function. *Proceedings of the National Academy of Sciences*, 107:4734-4739.

Boyacioğlu, R., Schulz, J., Koopmans, P. J., Barth, M., Norris, D. G. (2015). Improved sensitivity and specificity for resting state and task fMRI with multiband multi-echo EPI compared to multi-echo EPI at 7 T. *Neuroimage*, 119:352-361.

Bullmore, E., Sporns, O. (2009). Complex brain networks: graph theoretical analysis of structural and functional systems. *Nature Reviews Neuroscience*, 10:186-198.

Calhoun, V.D., Adali, T., Pearlson, G.D., Pekar, J.J. (2001). A method for making group inferences from functional MRI data using independent component analysis. *Human Brain Mapping*, 14:140–151.

Chand, G. B., Wu, J., Hajjar, I., Qiu, D. (2017). Interactions of the Salience Network and Its

Subsystems with the Default-Mode and the Central-Executive Networks in Normal Aging and Mild Cognitive Impairment. *Brain Connect*, 7:401-412.

Cohen, A. D., Yang, B., Fernandez, B., Banerjee, S., Wang, Y. (2020). Improved resting state functional connectivity sensitivity and reproducibility using a multiband multi-echo acquisition. *Neuroimage*, 225:117461.

D'Amour, A., Heller, K., Moldovan, D., et al. (2020). Underspecification presents challenges for credibility in modern machine learning. *arXiv*:2011.03395v2

Dong, Q., Ge, F., Ning, Q., Zhao, Y., Lv, J., Huang, H., Yuan, J., Jiang, X., Shen, D., Liu, T. (2020). Modeling Hierarchical Brain Networks via Volumetric Sparse Deep Belief Network. IEEE Trans Biomed Eng, 67:1739-1748.

Dummit, D. S., & Foote, R. M. (2004). *Abstract algebra* (Vol. 3). Hoboken: Wiley.

Duncan, J. (2010). The multiple-demand (MD) system of the primate brain: mental programs for intelligent behavior. *Trends in Cognitive Sciences*, 14:172-179.

Elliott, M. L., Romer, A., Knodt, A. R., Hariri, A. R. (2018). A connectome-wide functional signature of transdiagnostic risk for mental illness. *Biol Psychiatry*, 84:452-459.

Esteva, A., Robicquet, A., Ramsundar, B., Kuleshov, V., DePristo, M., Chou, K., Cui, C., Corrado, G., Thrun, S., & Dean, J. (2019).  A guide to deep learning in healthcare. *Nature Medicine*, 25:24-29.

Gurovich, Y., Hanani, Y., Bar, O., Nadav, G., Fleischer, N., Gelbman, D., ... & Bird, L. M. (2019). Identifying facial phenotypes of genetic disorders using deep learning. *Nature Medicine*, 25:60.

Hannun, A. Y., Rajpurkar, P., Haghpanahi, M., Tison, G. H., Bourn, C., Turakhia, M. P., & Ng, A. Y. (2019). Cardiologist-level arrhythmia detection and classification in ambulatory electrocardiograms using a deep neural network. *Nature Medicine*, 25(1), 65.

Hinton, G.E., Salakhutdinov, R.R. (2006). Reducing the dimensionality of data with neural networks. *Science*, 313:504-507.

Hinton, G.E., Osindero, S., Teh, Y.-W. (2006). A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527-1554.

Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N. (2012). Deep neural networks for acoustic modeling in

speech recognition: The shared views of four research groups. *Signal Processing Magazine, IEEE*, 29:82-97.

Hu, X., Huang, H., Peng, B., Han, J., Liu, N., Lv, J., ... & Liu, T. (2018). Latent source mining in FMRI via restricted Boltzmann machine. *Human Brain Mapping*, 39:2368-2380.

Huang, H., Hu, X., Zhao, Y., Makkie, M., Dong, Q., Zhao, S., ... & Liu, T. (2018). Modeling task fMRI data via deep convolutional autoencoder. *IEEE Transactions on Medical Imaging*, 37(7).

Hyvarinen, A. (1999). Fast and robust fixed-point algorithms for independent component analysis. *IEEE transactions on Neural Networks*, 10:626-634.

Iraji, A., Fu, Z., Damaraju, E., et al. (2019). Spatial dynamics within and between brain functional domains: A hierarchical approach to study time-varying brain function. *Hum Brain Mapp*, 40:1969-1986.

Iraji, A., Faghiri, A., Lewis, N., Fu, Z., Rachakonda, S., Calhoun, V. D. (2020). Tools of the trade: Estimating time-varying connectivity patterns from fMRI data. *Soc Cogn Affect Neurosci*, nsaa114. doi: 10.1093/scan/nsaa114.

Jilka, S.R., Scott, G., Ham, T., Pickering, A., Bonnelle, V., Braga, R. M., Leech, R., Sharp, D.J. (2014). Damage to the Salience Network and interactions with the Default Mode Network. *J Neurosci*, 34:10798-107807.

Kadison, R. V., & Ringrose, J. R. (1997). *Fundamentals of the theory of operator algebras* (Vol. 2). American Mathematical Soc..

Kaufmann, T., Alnæs, D., Doan, N. T,, Brandt, C. L., Andreassen, O. A,, Westlye, L.T. (2017). Delayed stabilization and individualization in connectome development are related to psychiatric disorders. *Nat Neurosci*, 20:513-515.

Kebets, V., Holmes A. J., Orban, C., Tang, S., Li, J., Sun, N., Kong, R., Poldrack, R. A., Yeo, B. T. T. (2019). Somatosensory-motor dysconnectivity spans multiple transdiagnostic dimensions of psychopathology. *Biol Psychiatry*, 86:779-791.

LeCun, Y., Bengio, Y., Hinton, G.E. (2015). Deep learning. *Nature*, 521:436-444.

Lee, D. D., & Seung, H. S. (1999). Learning the parts of objects by non-negative matrix factorization. *Nature*, 401:788-791.

Lee, K., Tak, S., Ye, J.C. (2011). A data-driven sparse GLM for fMRI analysis using sparse

dictionary learning with MDL criterion. *IEEE Transactions on Medical Imaging*, 30:1076-1089.

Lee, Y.-B., Lee, J., Tak, S., Lee, K., Na, D.L., Seo, S.W., Jeong, Y., Ye, J.C., Initiative, A.s.D.N. (2016). Sparse SPM: Group Sparse-dictionary learning in SPM framework for resting-state functional connectivity MRI analysis. *Neuroimage*, 125:1032-1045.

Li, L. M., Violante, I. R., Zimmerman, K., Leech, R., Hampshire, A., Patel, M., Opitz, A., McArthur, D., Jolly, A., Carmichael, D. W., Sharp, D. J. (2019). Traumatic axonal injury influences the cognitive effect of non-invasive brain stimulation. *Brain*,142:3280-3293.

Liu, J., Yuan, L., & Ye, J. (2010). An efficient algorithm for a class of fused lasso problems. *In Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining* (pp. 323-332). ACM.

Lv, J., Jiang, X., Li, X., Zhu, D., Zhang, S., Zhao, S., Chen, H., Zhang, T., Hu, X., Han, J. (2015). Holistic atlases of functional networks and interactions reveal reciprocal organizational architecture of cortical function. *Biomedical Engineering, IEEE Transactions on*, 62:1120-1131.

Mairal, J., Bach, F., Ponce, J., & Sapiro, G. (2010). Online learning for matrix factorization and sparse coding. *Journal of Machine Learning Research*, 11(1).

Mckeown, M. J., Sejnowski, T. J. (1998). Independent component analysis of fMRI data: Examining the assumptions. *Human Brain Mapping*, 6:368-372.

Menon, V., Toga, A. (2015). *Salience Network*. Elsevier. pp. 597–611. ISBN 978-0-12-397316-0.

Palacios, E. M., Yuh, E. L., Chang, Y. S., Yue, J. K., Schnyer, D. M., Okonkwo, D. O., Valadka, A. B., Gordon, W. A., Maas, A. I. R., Vassar, M., Manley, G. T., Mukherjee, P. (2017). Resting-State Functional Connectivity Alterations Associated with Six-Month Outcomes in Mild Traumatic Brain Injury. *J Neurotrauma*, 34:1546-1557.

Parkes, L., Satterthwaite, T. D,, Bassett, D. S. (2020). Towards precise resting-state fMRI biomarkers in psychiatry: synthesizing developments in transdiagnostic research, dimensional models of psychopathology, and normative neurodevelopment. *Curr Opin Neurobiol*, 65:120-128.

Plis, S. M., Hjelm, D. R., Salakhutdinov, R., Allen, E. A., Bockholt, H. J., Long, J. D., ... &

Calhoun, V. D. (2014). Deep learning for neuroimaging: a validation study. *Frontiers in Neuroscience*, 8:229.

Power, J. D., Cohen, A. L., Nelson, S. M., Wig, G. S., Barnes, K. A., Church, J. A., … & Petersen, S. E. (2011). Functional network organization of the human brain. *Neuron*, 72:665-678.

Royden, H. L. (1968). *Real analysis*. Krishna Prakashan Media.

Rudin, W. (1973). *Functional analysis*.

Schmidhuber, J. (2015). Deep learning in neural networks: An overview. Neural networks, 61:85-117.

Seo, J. D. (2018). "Deep" Independent Component Analysis in Tensorflow. https://towardsdatascience.com/deep-independent-component-analysis-in-tensorflow-manual-back-prop-in-tf-94602a08b13f

Shen, Y., Wen, Z., & Zhang, Y. (2014). Augmented Lagrangian alternating direction method for matrix separation based on low-rank factorization. *Optimization Methods and Software*, 29:239-263.

Smith, S. M., Fox, P. T., Miller, K. L., Glahn, D. C., Fox, P. M., Mackay, C. E., ... & Beckmann, C. F. (2009). Correspondence of the brain's functional architecture during activation and rest. *Proceedings of the National Academy of Sciences*, 106:13040-13045.

Sporns, O., Chialvo, D. R., Kaiser, M., & Hilgetag, C. C. (2004). Organization, development and function of complex brain networks. *Trends in Cognitive Sciences*, 8(9), 418-425.

Stam, C. J. (2014). Modern network science of neurological disorders. *Nature Reviews Neuroscience*, 15:683.

Suk, H.-I., Lee, S.-W., Shen, D., Initiative, A.s.D.N. (2014). Hierarchical feature representation and multimodal fusion with deep learning for AD/MCI diagnosis. *Neuroimage*, 101:569-582.

Suk, H.-I., Wee, C.-Y., Lee, S.-W., Shen, D. (2016). State-space model with deep learning for functional dynamics estimation in resting-state fMRI. *Neuroimage*, 129:292-307.

Trigeorgis, G., Bousmalis, K., Zafeiriou, S., & Schuller, B. W. (2016). A deep matrix factorization method for learning attribute representations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39:417-429.

Topol, E. J. (2019). High-performance medicine: the convergence of human and artificial intelligence. *Nature Medicine*, 25:44.

Vossel, S., Geng, J. J., Fink, G. R. (2014). Dorsal and ventral attention systems: distinct neural circuits but collaborative roles. *Neuroscientist*, 20:150-159.

Vu, A. T., Beckett, A., Setsompop, K., Feinberg, D. A. (2018). Evaluation of SLIce Dithered Enhanced Resolution Simultaneous MultiSlice (SLIDER-SMS) for human fMRI. *Neuroimage*, 164:164-171.

Wen, Z., Yin, W., & Zhang, Y. (2012). Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm. *Mathematical Programming Computation*, 4:333-361.

Zhang, W., Jiang, X., Zhang, S., Howell, B. R., Zhao, Y., Zhang, T., ... & Liu, T. (2017). Connectome-scale functional intrinsic connectivity networks in macaques. *Neuroscience*, *364*:1-14.

Zhang, W., Lv, J., Zhang, S., Zhao, Y., & Liu, T. (2018). Modeling resting state fMRI data via longitudinal supervised stochastic coordinate coding. *In Biomedical Imaging (ISBI 2018), IEEE 15th International Symposium on* (pp. 127-131). IEEE.

Zhang, W., Lv, J., Li, X., Zhu, D., Jiang, X., Zhang, S., ... & Liu, T. (2019). Experimental Comparisons of Sparse Dictionary Learning and Independent Component Analysis for Brain Network Inference from fMRI Data, *IEEE Transactions on Biomedical Engineering*, 66:289-299.

Zhang, W., Zhao, S., Hu, X., Dong, Q., Huang, H., Zhang, S., ... & Liu, T. (2020). Hierarchical Organization of Functional Brain Networks Revealed by Hybrid Spatiotemporal Deep Learning. *Brain Connectivity*, 10:72-82.

# Appendix A

1

2 *Assumption 1.1* For any operator discussed in this study, we have: $\forall \mathcal{C} \in \mathfrak{C}$, $\mathcal{C}: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$. This

3 assumption demonstrates that all operators are mapping from the finite dimensional space to another

4 finite dimensional space, which is also reasonable in the real world.

5 *Lemma 1.1* (**Norm Equality**) Given any arbitrary norm $\|\cdot\|$ and/or their finite linear combination

6 $\sum_{i=1}^{n} k_i \|\cdot\|$ denoted based on any finite set, this norm or their finite linear combination is equivalent to

7 $\ell_2$ norm (e.g., $\|\cdot\|_2$).

8 *Proof*: We denote $\ell_1$ and $\ell_2$ norm in finite dimensions, such as:

$$\|a\|_1 = \sum_{i=1}^{n} |a_i|$$

$$\|a\|_2 = \left(\sum_{i=1}^{n} a_i^2\right)^{\frac{1}{2}} \quad \text{(A.1)}$$

$$a = [a_1, a_2, \cdots, a_n]$$

9 Obviously, since all norms are non-negative, according to Eq. (A.1), we have:

$$\sum_{i=1}^{n} a_i^2 \leq \left(\sum_{i=1}^{n} |a_i|\right)^2 \quad \text{(A.2)}$$

10 Eq. (A.2) implies:

$$\|a\|_2 \leq \|a\|_1 \quad \text{(A.3)}$$

11 And, based on Cauchy-Schwarz inequality, we have:

$$\|a\|_1^2 = \left(\sum_{i=1}^{n} |a_i| \cdot 1\right)^2 \leq \sum_{i=1}^{n} a_i^2 \cdot \sum_{j=1}^{n} 1^2 = \|a\|_2^2 \cdot n \quad \text{(A.4)}$$

12 It implies:

$$\frac{1}{\sqrt{n}} \|a\|_1 \leq \|a\|_2 \quad \text{(A.5)}$$

13 According to the theorem of norm equality (Rudin, 1973), given an arbitrary finite dimensional space,

14 if and only if the following inequality holds:

$$c\|\cdot\|_2 \leq \|\cdot\| \leq C\|\cdot\|_2 \quad \text{(A.6)}$$

15 Thus, the norm $\|\cdot\|$ is equivalent to $\|\cdot\|_2$. Since Eq. (B.3) and Eq. (B.4) hold, we have:

$$c\frac{1}{\sqrt{n}} \|a\|_1 \leq \|a\|_2 \leq \|a\|_1 \quad \text{(A.7)}$$

40

1  It implies $\|\cdot\|_1$ and $\|\cdot\|_2$ are equivalent. Similarly, we can prove $\sum_{i=1}^{n} k_i \|\cdot\|$ is also equivalent to $\|\cdot\|_2$.

2

3  ***Theorem 1.1* (Superiority of Deep Linear Models)** Given a real function $f(x)$ and $m(\{x \in$

4  $[a,b]: |f(x)| = \pm\infty\}) = 0$. If considering the series of polynomials $\{P_n(x)\}_{n=1}^{N}$, we have: if $N$ is

5  large enough, we have: $\forall \varepsilon > 0$ $\left\|\{P_n(x)\}_{n=1}^{N} - f(x)\right\| \le \varepsilon$; if $N \to \infty$, we have: $\lim_{N\to\infty}\{P_n(x)\}_{n=1}^{N} =$

6  $f(x)$; however, for any shallow model, since $N$ should be bounded, we only have: $\left\|\{P_n(x)\}_{n=1}^{N} -$

7  $f(x)\right\| \le M$.

8  ***Proof***: According to Лузин (Luzin) Theorem (Royden, 1968), we have a close set:

$$F_n \subset F_{n+1} \subset \cdots \subseteq [a,b]$$
$$m([a,b]\backslash F_n) = \frac{1}{n} \tag{A.8}$$
$$f \in C(F_n)$$

9  Then we have a consistent real function $g(x)$, and obviously we have:

$$g(x) = f(x) \tag{A.9}$$

10  Since for any continuous real function, we have:

$$|g(x) - P_n(x)| < \frac{1}{n} \tag{A.10}$$

11  Let $\mathcal{F} = \bigcup_{n=1}^{\infty} F_n$, and obviously we have:

$$m([a,b]\backslash \mathcal{F}) = 0 \tag{A.11}$$

12  If $\mathfrak{F}$ is a real function denoted on set $\mathcal{F}$, it indicates:

$$\lim_{N\to\infty}\left|\mathfrak{F} - \{P_n(x)\}_{n=1}^{N}\right| < \varepsilon \tag{A.12}$$

13  then we have $\lim_{N\to\infty}\{P_n(x)\}_{n=1}^{N} = \mathfrak{F}$,

14  Moreover, it is easy to prove $\{P_n\}_{n=1}^{\infty}$ denoted on $[a,b]$ as a ring $(\{P_n\}_{n=1}^{\infty}, +, \times)$ (Dummit, 2004;

15  Kadison, 1997). And $m(\cdot)$ represents a Lebesgue measure.

16

17  ***Lemma 1.2* (Contraction of Operators Combination)** Given two contraction mappings $\Phi_1$ and $\Phi_2$,

18  we have the composite of two contraction mapping as $\Phi_2 \cdot \Phi_1$. The composite mapping $\Phi_2 \cdot \Phi_1$ must

19  be contractive.

20  ***Proof***: According to the definition of contraction linear operator, we have:

$$\exists \zeta \in (0,1)$$
$$\rho \stackrel{\text{def}}{=} \|\Phi x - \Phi y\|$$
$$\rho(\Phi x, \Phi y) \le \zeta \rho(x,y) \tag{A.13}$$

1     Obviously, and we have:

$$\rho(\Phi_1 u, \Phi_1 v) \leq \zeta \rho(u,v) \ \ \forall \zeta \in (0,1)$$
$$\rho(\Phi_2 x, \Phi_2 y) \leq \eta \rho(x,y) \ \ \forall \eta \in (0,1)$$

(A.14)

2     If we set:

$$x = \Phi_1 u, y = \Phi_1 v$$

(A.15)

3     the inequality below holds:

$$\rho(\Phi_2 x, \Phi_2 y) \leq \eta \rho(\Phi_1 u, \Phi_1 v) \leq \zeta \eta \rho(u,v)$$

(A.16)

4     Since the definition as

$$\forall \zeta, \eta \in (0,1), \ \rho(\Phi_2 \Phi_1 u, \Phi_2 \Phi_1 y) \leq \zeta \eta \rho(u,v)$$

(A.17)

5

6     ***Corollary 1.1*** **(General Contraction Operator)** According to Lemma 1.2, if denote the operators
7     $\{\Phi_i\}_{i=1}^K$, $\forall \Phi_i \ i \in \mathbb{N}$, $\Phi_i: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$; considering any combination of operators: $\Phi_K \cdot \cdots \cdot \Phi_2 \cdot \Phi_1$,
8     if at least a single operator $\Phi_i$ is contraction operator, and other operators are bounded, such as $\forall i \neq$
9     $k \ \|\Phi_i\| \leq M$. If and only if $\prod_{i=1}^K \|\Phi_i\| < 1$, the combination of operator series $\Phi_K \cdot \cdots \cdot \Phi_2 \cdot \Phi_1$ is a
10     contraction operator.

11     ***Proof***: Obviously, according to Lemma 1.2, use a series as $\{\zeta_i\}_{i=1}^K$ to replace $\zeta, \eta \in (0,1)$,

12     Obviously, we have:

$$\zeta_i \in (0,1) \ i \in \mathbb{N}$$
$$\rho(\Phi_K \cdot \cdots \cdot \Phi_2 \Phi_1 u, \Phi_K \cdot \cdots \cdot \Phi_2 \Phi_1 y) \leq \zeta_K \cdot \cdots \zeta_2 \cdot \zeta_1 \cdot \rho(u,v)$$

(A.18)

13

14     Since $\zeta_K \cdot \cdots \zeta_2 \cdot \zeta_1 < 1$, we have proved this corollary.

15

16     ***Corollary 1.2*** **(Iterative Contraction Operator)** According to Lemma 1.2, if denote the operators
17     $\{\Phi_i\}_{i=1}^K$, $\forall \Phi_i \ i \in \mathbb{N}$, $\Phi_i: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$; considering any combination of operators: $\Phi_K \cdot \cdots \cdot \Phi_2 \cdot \Phi_1$,
18     if at least a single operator $\Phi_i$ is contraction operator, and other operators are bounded, such as $\forall i \neq$
19     $k, \|\Phi_i\| \leq M$. If and only if $\lim_{n \to \infty} \prod_{i=1}^K \|\Phi_i\|^n = c < 1$, the combination of operator series $\Phi_K^n \cdot \cdots \cdot \Phi_2^n \cdot$
20     $\Phi_1^n$.

21     ***Proof***: Obviously, according to Lemma 1.2 and Corollary 1.1 and 1.2, use a series as $\{\zeta_i\}_{i=1}^K$ to replace
22     $\zeta, \eta \in (0,1)$,

23     And we have:

$$\forall \zeta_i \in (0,1) \ i \in \mathbb{N}$$
$$\rho(\Phi_K^n \cdot \cdots \cdot \Phi_2^n \cdot \Phi_1^n u, \Phi_K^n \cdot \cdots \cdot \Phi_2^n \cdot \Phi_1^n y) < \zeta_i^n \cdot \cdots \cdot \zeta_2^n \cdot \zeta_1^n \cdot \rho(u,v)$$

(A.19)

24

25     Since $0 < \zeta_i^n \cdot \cdots \cdot \zeta_2^n \cdot \zeta_1^n < 1$, we have proved this corollary.

26

# Appendix B

1

2

3    *Definition 2.1* If we denote Deep MF as an operator $\mathfrak{A}$, based on the description of Deep MF,

4    considering the iteration k, we can denote $\mathfrak{A} \overset{\text{def}}{=} M \cdot \mathcal{A}^{\text{k}} \cdot \mathcal{S}^{\text{k}} \cdot \mathcal{R}^{\text{k}}$.

5    *Definition 2.2* If we denote Deep SDL as an operator $\mathfrak{L}$, based on the description of Deep SDL,

6    considering the iteration k, we can denote $\mathfrak{L} \overset{\text{def}}{=} M \cdot \mathcal{G}^{\text{k}} \cdot \mathcal{S}^{\text{k}}$.

7    *Definition 2.3* If we denote Deep FICA as an operator $\mathfrak{T}$, based on the description of Deep FICA,

8    considering the iteration k, we can denote $\mathfrak{T} \overset{\text{def}}{=} \mathcal{P} \cdot \mathcal{F}^{\text{k}}$.

9    *Definition 2.4* If we denote Deep NMF as an operator $\mathfrak{N}$, based on the description of Deep NMF,

10    considering the iteration k, we can denote $\mathfrak{N} \overset{\text{def}}{=} M \cdot \mathcal{U}^{\text{k}} \cdot \mathcal{N}$.

11

12    *Theorem 2.1* (**Contraction of ADMM Operator**) ADMM could be considered as contraction operator.

13    It can be treated as a general iterative contraction operator in finite dimensionality space. We have

14    ADMM $\overset{\text{def}}{=} \mathcal{A}$. If denote the $\left\| \mathcal{A}^{k+1} \right\| = \alpha \left\| \mathcal{A}^{k} \right\|$, and $\beta$ should be step length, i.e., penalty parameter,

15    if $n \to \infty \;\; 0 < (\alpha\beta)^{n} \|BN\| < 1$ , $\mathcal{A}$ can be considered as a contraction operator. And $\|BN\|$ denotes

16    the norm of different residual error, considering two distinctive input matrices.

17    *Proof:* $X$ and $Y$, represent the two input matrices.

18    Consider the iterative format of ADMM as

$$\mathcal{A}_{k+1} \leftarrow \mathcal{A}_{k} - \min(f_{\mathcal{A}}) \tag{B.1}$$

19    And it also can imply:

$$\begin{aligned} \|\mathcal{A}_{k+1}\| = \alpha\|\mathcal{A}_{k}\|, \\ 0 < \alpha < 1 \end{aligned} \tag{B.2}$$

20    According to the definition of contraction operator, we have:

$$\|\mathcal{A}X - \mathcal{A}Y\| \le \alpha \left\| \left( \beta \left( \mathbf{e}_{k}^{t} + \prod_{i=1}^{k-1} X_{i}\,Y_{k} + \sum_{i=1}^{k} Z_{k}^{t+1} - SG \right) - \alpha\beta \left( \hat{\mathbf{e}}_{k}^{t} + \prod_{i=1}^{k-1} \hat{X}_{i}\,\hat{Y}_{k} \right. \right. \\ \left. \left. + \sum_{i=1}^{k} \hat{Z}_{k}^{t+1} - \widehat{SG} \right) \right) \right\| \tag{B.3}$$

21    And we also have:

$$\|\mathcal{A}X - \mathcal{A}Y\| \le \alpha\beta \left\| \mathbf{e}_{k}^{t} - \hat{\mathbf{e}}_{k}^{t} + \prod_{i=1}^{k-1} X_{i}\,Y_{k} - \prod_{i=1}^{k-1} \hat{X}_{i}\,\hat{Y}_{k} + \sum_{i=1}^{k} Z_{k}^{t+1} - \sum_{i=1}^{k} \hat{Z}_{k}^{t+1} \right. \\ \left. + \widehat{SG} - SG \right\| \tag{B.4}$$

43

1   Since $\mathfrak{e}_k^t$, $\hat{\mathfrak{e}}_k^t$, $\prod_{i=1}^{k-1} X_i Y_k$, $\prod_{i=1}^{k-1} \hat{X}_i \hat{Y}_k$, $\sum_{i=1}^{k} Z_k^{t+1}$, $\sum_{i=1}^{k} \hat{Z}_k^{t+1}$, $\hat{S}, S \in \mathbb{R}^{m \times n}$ , they are obviously
2   bounded; and using Corollary 1.1 and 1.2, we have:

$$\|\mathcal{A}X - \mathcal{A}Y\| \leq \alpha\beta\|BN\| \tag{B.5}$$

3   Obviously, it demonstrates:

$$\|\mathcal{A}^n A - \mathcal{A}^n B\| \leq (\alpha\beta)^n\|BN\| < 1 \tag{B.6}$$

4   If and only if $0 < (\alpha\beta)^n < 1$, or $0 < \alpha\beta < 1$, $\mathfrak{A}^n$ is equivalent to a contraction operator. According
5   to Lemma 1.2 and Corollary 1.1, 1.2, it also indicates: when $n$ is large enough, $n > N$, we have:

$$\lim_{n\to\infty}\|\mathcal{A}^n A - \mathcal{A}^n B\| \leq \lim_{n\to\infty}(\alpha\beta)^n\|BN\| \tag{B.7}$$

6   Obviously, if and only if $\lim_{n\to\infty}(\alpha\beta)^n\|BN\| < 1$, the iterative ADMM operator can be equivalent to a

7   contraction operator.

8

9   ***Theorem 2.2 (Initialization Operator is bounded)*** If we denote the sparse operator as $\mathcal{M}: \mathbb{R}^{S \times T} \to$
10   $\mathbb{R}^{S \times T}$, we have $\|\mathcal{M}\| < \infty$.

11   ***Proof***: according to the definition of operator norm (Rudin 1973), $\|\mathcal{M}\| \leq sup\frac{\|MX\|}{\|X\|}$; obviously,

12   $\|\mathcal{M}X\|$ and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we
13   denote:

$$X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} \quad \mathcal{M} = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix} \quad \|X\| < \infty \quad \|\mathcal{M}X\| < \infty \tag{B.8}$$

14   Obviously, $\|\mathcal{M}\| < \infty$.

15

16   ***Theorem 2.3 (Sparsity Operator is bounded)*** If we denote the sparse operator as $\mathcal{S}: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$,
17   we have $\|\mathcal{S}\| < \infty$.

18   ***Proof***: according to the definition of operator norm (Rudin, 1973), $\|\mathcal{S}\| \leq sup\frac{\|SX\|}{\|X\|}$; obviously, $\|\mathcal{S}X\|$

19   and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we denote:

$$\mathcal{S}X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ 0 \\ a_n \end{bmatrix} \quad \mathcal{S}Y = \begin{bmatrix} b_1 \\ 0 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix} \tag{B.9}$$

20   and we examine:

$$\mathcal{S}X - \mathcal{S}Y = \begin{bmatrix} a_1 - b_1 \\ a_2 \\ \vdots \\ - b_{n-1} \\ a_n - b_n \end{bmatrix}; \quad X - Y = \begin{bmatrix} a_1 - b_1 \\ a_2 - b_2 \\ \vdots \\ a_{n-1} - b_{n-1} \\ a_n - b_n \end{bmatrix}, \quad \|\mathcal{S}X - \mathcal{S}Y\| \le s \ \|X - Y\|, \quad \text{(B.10)}$$

1   Without loss of generality, and based on Lemma 1.2, we calculate the $\ell_2$ norm, and we have:

$$\infty > s \ge \frac{\sum_{i=u}^{n}(a_i - b_i)^2 + \sum_{i=v}^{p}(a_i)^2 + \sum_{i=w}^{t}(b_i)^2}{\sum_{i=1}^{n}(a_i - b_i)^2} \tag{B.11}$$

2   This inequality demonstrates that $\|\mathcal{S}\|$ is a bounded. And $\mathcal{S}$ is a bounded operator.

3

4   ***Theorem 2.4* (Rank Reduction Operator is bounded)** If we denote the sparse operator as
5   $\mathcal{R}: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$, we have $\|\mathcal{R}\| < \infty$.

6   ***Proof***: According to the definition of operator norm (Rudin, 1973), $\|\mathcal{R}\| \le sup \frac{\|\mathcal{R}X\|}{\|X\|}$; obviously, $\|\mathcal{R}X\|$
7   and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we denote:

$$X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix}, \quad \mathcal{R}X = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_k \\ \vdots \\ 0 \end{bmatrix} \tag{B.12}$$

8   Eq. (B.31) implies:

$$sup \frac{\|\mathcal{R}X\|}{\|X\|} = \frac{\sum_{i=1}^{n} a_i^2}{\sum_{i=u}^{p}(a_i - b_i)^2 + \sum_{i=v}^{q} a_i^2} < \infty. \tag{B.13}$$

9

10   ***Theorem 2.6* (Normalization Operator of Deep NMF is bounded)** If we denote the normalization
11   operator of Deep NMF as $\mathcal{N}: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$, we have $\|\mathcal{N}\| \le 1$.

12   ***Proof***: according to the definition of operator norm (Rudin, 1973), $\|\mathcal{N}\| \le sup \frac{\|\mathcal{N}X\|}{\|X\|}$; obviously,
13   $\|\mathcal{N}X\|$ and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we
14   denote:

$$X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} \quad \mathcal{N}X = \begin{bmatrix} b_1 \\ 0 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix} \tag{B.14}$$

15

16   According to Eq. (B.33), we need to notice: $\{a_i\}_{i=1}^{K} \subseteq [-q, q]$, $1 \le q < \infty$; $\{b_i\}_{i=1}^{K} \subseteq [0,1]$.
17   Obviously, $\|\mathcal{N}X\| < \|X\|$. Finally, we have: $\|\mathcal{N}\| < 1$.

18

1    ***Theorem 2.7 (Contraction of Updating Operator Deep NMF)*** If we denote the updating operator as

2    $\mathcal{U}: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$, we have $\|\mathcal{U}\| < 1$.

3    ***Proof***: according to the definition of operator norm (Rudin, 1973), $\|\mathcal{U}\| \le \sup \frac{\|\mathcal{U}X\|}{\|X\|}$; obviously, $\|\mathcal{U}X\|$

4    and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we denote:

$$X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} \quad \mathcal{U}X = \begin{bmatrix} b_1 \\ 0 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix} \tag{B.15}$$

5    According to the iterative format of Deep NMF, we need to notice: $b_i = \frac{a_i}{\max f(a_i)}$; obviously, $\|\mathcal{U}X\| <$

6    $\|X\|$. Finally, we have: $\|\mathcal{U}\| < 1$. Otherwise, if $\|\mathcal{U}\| > 1$, when $k \to \infty$, we have: $\|\mathcal{U}X\| = \infty$.

7

8

9    ***Theorem 2.8 (Contraction of GD Operator)*** Gradient Descent (GD) is a bounded contraction

10   operator, if and only if the derivative of target function is bounded:

11   $|f''(\varsigma)| < \frac{1}{\sigma} < \infty$, $\sigma$ is the step length.

12   ***Proof***: The standard iteration format is:

$$x_{k+1} = x_k - \sigma f'(x_k) \tag{B.16}$$

13   Using the definition of operator, we have:

$$\tau(x_k) = x_k - \sigma f'(x_k) \quad \forall \sigma \in (0,1) \tag{B.17}$$

14   And we have:

$$\tau \|\tau X - \tau Y\| = \|(X - Y) - \sigma(f'(X) - f'(Y))\| \tag{B.18}$$

15   Using Mean value theorem, we have:

$$\tau \|\tau X - \tau Y\| = |1 - \sigma f''(\varsigma)| \|X - Y\| \tag{B.19}$$

16   According to the definition of contraction operator (Rudin, 1973), if and only if:

$$|1 - \sigma f''(\varsigma)| < 1, \ |1 - \sigma f''(\varsigma)| \in \mathbb{K} \tag{B.20}$$

17   It also implies, when the following inequality holds:

$$|f''(\varsigma)| < \frac{1}{\sigma} < \infty \tag{B.21}$$

18   GD is considered as a contraction mapping/operator. Without generality, we can set $\sigma < \frac{1}{|f''(x)|+1}$.

19   And obviously, using multiplicative inequality, we have:

46

$$\|\tau X - \tau Y\| \leq \|\tau\| \|X - Y\| \tag{B.22}$$

1.  Since $X$ and $Y$ both denote in finite $\ell^2$ space, we have:

$$\|\tau\| \|X - Y\| \leq \infty \tag{B.23}$$

2.  Using Uniformly bounded theorem, we have:

$$\|\tau\| \leq M, \ M \in \mathbb{K} \tag{B.24}$$

3.  GD is a bounded mapping/operator.

4.  According to Lemma 1.2, and Corollary 1.1-1.2, obviously, for $n$ iterations for an operator, and if we
5.  set the accuracy level as $\varepsilon$, we have:

$$\|\tau^n X - \tau^{n+1} Y\| = \sigma^n \|X - \tau Y\| < \varepsilon \tag{B.25}$$

6.  Since $X$ and $Y$ is both denoted in finite $\ell^2$ space, we have:

$$\sigma^n \|X - \tau Y\| \leq \sigma^n (\|X\| + \|\tau Y\|) \tag{B.26}$$

7.  Obviously, $\|X - Y\|_{\ell^2}$ is bounded, and we have:

$$\sigma^n(\|X\| + \|\tau Y\|) \leq \sigma^n(\|X\| + \|\tau\| \|Y\|) \leq \sigma^n(\|X\| + \|Y\|) \leq \sigma^n \cdot 2\|X\|$$
$$0 < \sigma^n \cdot 2\|X\| < \varepsilon$$
$$n > log \frac{\varepsilon}{2\|X\|} \ / \ log \ \sigma > 0 \tag{B.27}$$

8.  We provide the infimum of iteration as $log \frac{\varepsilon}{2\|X\|} \ / \ log \ \sigma$ to approach the accuracy level $\varepsilon$.

9.

10. ***Theorem 2.9* (Operator PCA is bounded)** If we denote the updating operator as $\mathcal{P}: \mathbb{R}^{S \times T} \to \mathbb{R}^{S \times T}$,
11. we have $\|\mathcal{P}\| < \infty$.

12. ***Proof***: According to the definition of operator norm (Rudin, 1973), $\|\mathcal{P}\| \leq sup \frac{\|\mathcal{P}X\|}{\|X\|}$; obviously, $\|\mathcal{U}X\|$

13. and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we denote:

$$X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} \quad \mathcal{P}X = \begin{bmatrix} b_1 \\ b_2 \\ \vdots \\ b_{n-k} \\ \vdots \\ 0 \end{bmatrix} \tag{B.29}$$

14. According to the dimensional reduction of PCA, we have: $sup \frac{\|\mathcal{P}X\|}{\|X\|} = sup \frac{(\sum_{i=1}^{n-k} b_i^2)^{\frac{1}{2}}}{(\sum_{i=1}^{n} a_i^2)^{\frac{1}{2}}} < \infty$. It

15. demonstrates: $\|\mathcal{P}\| < \infty$.

16.

1  **_Theorem 2.10_ (Contraction of Fixed-Point Operator)** If we denote the updating operator as

2  $\mathcal{F}: \mathbb{R}^{S \times T} \rightarrow \mathbb{R}^{S \times T}$, we have $\|\mathcal{F}\| < 1$.

3  **_Proof_**: according to the definition of operator norm (Rudin, 1973), $\|\mathcal{F}\| \leq sup \frac{\|\mathcal{F}X\|}{\|X\|}$; obviously, $\|\mathcal{U}X\|$

4  and $\|X\|$ is bounded, since both of norms are based on finite dimensional matrix. And if we denote:

$$X = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_{n-1} \\ a_n \end{bmatrix} \quad \mathcal{F}X = \begin{bmatrix} b_1 \\ 0 \\ \vdots \\ b_{n-1} \\ b_n \end{bmatrix} \tag{B.30}$$

5  According to the iterative format of Deep NMF, we need to notice: $b_i = \frac{a_i}{\sqrt{\|a_i\| C_i \|a_i^T\|}}$; obviously,

6  $\|\mathcal{F}X\| < \|X\|$. Finally, we have: $\|\mathcal{F}\| < 1$. Otherwise, if $\|\mathcal{F}\| > 1$, when $k \rightarrow \infty$, we have: $\|\mathcal{F}X\| = $

7  $\infty$.

8

9

10  **_Theorem 2.11_ (Inequality of Operator Norms)** According to Theorem 2.1, 2.7, 2.8 and 3.0, if we

11  assume: $\|\mathcal{A}^{k+1}\| = \alpha_1 \|\mathcal{A}^k\|$, $\|\mathcal{G}^{k+1}\| = \alpha_2 \|\mathcal{G}^k\|$, $\|\mathcal{U}^{k+1}\| = \alpha_3 \|\mathcal{U}^k\|$, $\|\mathcal{F}^{k+1}\| = \alpha_4 \|\mathcal{F}^k\|$, we

12  have: $\alpha_1 \neq \alpha_3, \alpha_4$; $\alpha_2 \neq \alpha_3, \alpha_4$;

13  **_Proof_**: Proof by contradiction. In general, we assume $\|\mathcal{A}\| = \|\mathcal{U}\|$, according to the iterative formats

14  of Deep MF and Deep NMF, and considering:

$$\mathcal{A}_{k+1} \leftarrow \mathcal{A}_k - \min(f_{\mathcal{A}}) \tag{B.31}$$

15  If we employ the $\alpha\mathcal{A}_k = \mathcal{A}_{k+1}$ to replace $\mathcal{A}_{k+1}$:

$$\alpha\mathcal{A}_k = \mathcal{A}_k - \min(f_{\mathcal{A}}) \tag{B.32}$$

16  And we can reformat this equality as:

$$(1 - \alpha)\mathcal{A}_k = \min(f_{\mathcal{A}}) \tag{B.33}$$

17  Considering the iterative format of Deep NMF:

$$\mathfrak{N}_{k+1} \leftarrow \mathfrak{N}_k / max(f_{\mathfrak{N}}) \tag{B.34}$$

18  Let we denote:

$$\left| \frac{1}{max(f_{\mathfrak{N}})} \right| \leq \varepsilon \tag{B.35}$$

19  And considering an extreme condition, $\forall \varepsilon_i \leq \varepsilon$, for each iteration $i$, and $\lim_{i \to \infty} \varepsilon_i = \varepsilon$;

$$\lim_{i \to \infty}(1 - \varepsilon_i)\mathcal{A}_k = \min(f_\mathcal{A}) \tag{B.36}$$

1

2    Then we have the conclusion:

$$\exists\ n \ll k, \mathcal{A}_k = \min(f_\mathcal{A}) \tag{B.37}$$

3    It demonstrates for the iterative format of Deep MF, before convergence, the iteration can be terminated,

4    since a very small norm of operator $\mathcal{A}$. $\mathcal{A}$ cannot guarantee the convergence. It obviously disobeys

5    the property of ADMM.

6    Similarly, we can also prove $\alpha_2 \neq \alpha_3$, $\alpha_4$; and $\alpha_1 \neq \alpha_4$.

7

1 # Appendix C

2

3 ***Assumption 3.1*** For all operators, these operators should be considered as linear operators, and we have:

$$\Phi \cdot (X + Y) = \Phi \cdot X + \Phi \cdot Y \tag{C.1}$$

4 ***Assumption 3.2*** For any input matrix, we can successfully separate the vital information and

5 background noise. If we denote: $V = \{\bigcup_{i=1}^{P} voxel_i, \ voxel_i \in BN\}$, and $N = \{\bigcup_{i=1}^{Q} voxel_i, \ voxel_i \notin$

6 $BN\}$. BN represents the functional areas, i.e., potentially activated areas of brain. We have some crucial

7 assumptions: $V \cap N = \emptyset, \ V \succcurlyeq 0, \ B \succcurlyeq 0, \ \|V\| \gg \|N\|$.

8 ***Lemma 3.1*** (**Continuous Operators**) For all operators analyzed in this study, if $k > K, \forall k \in \mathbb{N}$, these
9 iterative operators can be considered as consistent operator. It means: if we have $\|V - \hat{V}\| \le \varepsilon$,
10 $\|\mathfrak{A}^k V - \mathfrak{A}^k \hat{V}\| \to 0$.

11 ***Proof***: We denote: $\mathfrak{A}, \mathfrak{L}, \mathfrak{T}, \mathfrak{N} \in \mathfrak{C}$: $\mathbb{R}^{s \times t} \to \mathbb{R}^{s \times t}$

12 For $V, \hat{V} \subseteq \mathbb{R}^{s \times t}$, we assume that:

$$\|V - \hat{V}\| \le \frac{\varepsilon}{M} \tag{C.2}$$

13 If $k > K$, For any operator belongs to $\mathfrak{C}$ can be considered as a contraction operator, and we have:

$$\|\mathfrak{A}^k V - \mathfrak{A}^k \hat{V}\| \le \|\mathfrak{A}^k\| \cdot \|V - \hat{V}\| \le M \cdot \frac{\varepsilon}{M} = \varepsilon \tag{C.3}$$

14 This inequality demonstrates that all operators of $\mathfrak{C}$, if $k$ is large enough, can be treated as the consistent
15 operators (Rudin, 1973). Similarly, it also demonstrates: $\|\mathfrak{A}^k V - \mathfrak{L}^k V\| \le \varepsilon$

16

17 ***Theorem 3.1*** (**Distinctive Spatial Similarity**) If we denote the following set:

$$Deep \ MF: D = \{\mathfrak{A}^k N, N \in \bigcup_{i=1}^{M} voxel_i, \ voxel_i \notin T\}$$

$$Deep \ SDL: L = \{\mathfrak{L}^k N, N \in \bigcup_{i=1}^{M} voxel_i, \ voxel_i \notin T\}$$

$$Deep \ FICA: I = \{\mathfrak{T}^k N, N \in \bigcup_{i=1}^{M} voxel_i, \ voxel_i \notin T\} \tag{C.4}$$

$$Deep \ NMF: \Theta = \{\mathfrak{N}^k N, N \in \bigcup_{i=1}^{M} voxel_i, \ voxel_i \to 0\}$$

18 And considering the iteration *k*, it implies:

$$\frac{|\mathfrak{A}^k V|}{|V \cup D|} \le \frac{|\mathfrak{L}^k V|}{|V \cup L|} \le \frac{|\mathfrak{T}^k V|}{|V \cup I|} \le \frac{|\mathfrak{N}^k V|}{|V \cup \Theta|} \tag{C.5}$$

1    where $|\cdot|$ denotes the number of positive elements.

2    ***Proof***:  Based on assumptions 3.1 and 3.2, if $\forall\ k \in \mathbb{N}$, we have:

$$\begin{aligned}
\mathfrak{A}^k C &= \mathfrak{A}^k V + (\mathfrak{A}^k N) \\
\mathfrak{L}^k C &= \mathfrak{L}^k V + (\mathfrak{L}^k N) \\
\mathfrak{T}^k C &= \mathfrak{T}^k V + (\mathfrak{T}^k N) \\
\mathfrak{N}^k C &= \mathfrak{N}^k V + \Theta
\end{aligned} \tag{C.6}$$

3    According to Corollary 1.1 and 1.2, $k > K$, we have:

$$0 = \|\Theta\| \le \|\mathfrak{T}^k N\| \le \|\mathfrak{L}^k N\| \le \|\mathfrak{A}^k N\| < \infty \tag{C.7}$$

4    We can also rewrite it as:

$$0 = |\Theta| < |I| \le |L| \le |D| \tag{C.8}$$

5    And, according to the spatial similarity, we also have:

$$\begin{aligned}
Deep\ MF_{Similarity} &\overset{\text{def}}{=} \frac{\left|(\mathfrak{A}^k V \cup D) \cap V\right|}{|V \cup (\mathfrak{A}^k A \cup D)|} = \frac{|\mathfrak{A}^k V|}{|V \cup D|} \\
Deep\ SDL_{Similarity} &\overset{\text{def}}{=} \frac{\left|(\mathfrak{L}^k V \cup L) \cap A\right|}{|V \cup (\mathfrak{L}^k V \cup L)|} = \frac{|\mathfrak{L}^k A|}{|V \cup L|} \\
Deep\ FICA_{Similarity} &\overset{\text{def}}{=} \frac{\left|(\mathfrak{T}^k V \cup I) \cap V\right|}{|V \cup (\mathfrak{T}^k A \cup I)|} = \frac{|\mathfrak{T}^k V|}{|V \cup I|} \\
Deep\ NMF_{Similarity} &\overset{\text{def}}{=} \frac{\left|(\mathfrak{N}^k V \cup \Theta) \cap V\right|}{|V \cup (\mathfrak{N}^k A \cup \Theta)|} = \frac{|\mathfrak{N}^k V|}{|V|}
\end{aligned} \tag{C.10}$$

6    Again, considering $k > K$, and Corollary 1.1 to 1.2, and Theorem 3.2, we have:

$$\left|\mathfrak{N}^k V\right| = \left|\mathfrak{T}^k V\right| = \left|\mathfrak{L}^k V\right| = \left|\mathfrak{A}^k V\right| \tag{C.11}$$

7    Obviously, we have:

$$0 < |V| = |V \cup \Theta| \le |V \cup I| \le |V \cup L| \le |V \cup D| < \infty \tag{C.12}$$

8    Finally, the following inequality holds, such that:

$$0 < \frac{|\mathfrak{A}^k V|}{|V \cup D|} \le \frac{|\mathfrak{L}^k V|}{|V \cup L|} \le \frac{|\mathfrak{T}^k V|}{|V \cup I|} \le \frac{|\mathfrak{N}^k V|}{|V \cup \Theta|} \tag{C.13}$$

9

10    ***Theorem 3.2*** (**Bounded Iterative Operators**) For all operators analyzed in this study, if $k > K, \forall k \in$
11    $\mathbb{N}$, these iterative operators can be considered as consistent operator. If we have: $\|\hat{V}\| \le \varepsilon$, it means:
12    $\|\mathfrak{N}^k V - \mathfrak{A}^k V\| \to 0$, $\|\mathfrak{L}^k V - \mathfrak{A}^k V\| \to 0$ and $\|\mathfrak{T}^k V - \mathfrak{A}^k V\| \to 0$.

13    ***Proof***:  We denote: $\mathfrak{A}, \mathfrak{L}, \mathfrak{T}, \mathfrak{N} \in \mathfrak{C}: \mathbb{R}^{s \times t} \to \mathbb{R}^{s \times t}$

14    If $k > K$, For any operator belongs to $\mathfrak{C}$ can be considered as a contraction operator, according to
15    Lemma 3.1, and we have:

$$\left\| \mathfrak{N}^k V - \mathfrak{A}^k V \right\| = \left\| \mathfrak{N}^k V - \mathfrak{N}^k \hat{V} + \mathfrak{N}^k \hat{V} - \mathfrak{A}^k \hat{V} + \mathfrak{A}^k \hat{V} - \mathfrak{A}^k V \right\| \tag{C.14}$$

$$\leq \left\| \mathfrak{N}^k V - \mathfrak{N}^k \hat{V} \right\| + \left\| \mathfrak{N}^k \hat{V} - \mathfrak{A}^k \hat{V} \right\| + \left\| \mathfrak{A}^k \hat{V} - \mathfrak{A}^k V \right\|$$

1   According to Lemma 3.1, and we have:

$$\left\| \mathfrak{N}^k V - \mathfrak{N}^k \hat{V} \right\| \leq \frac{\varepsilon}{3}$$

$$\left\| \mathfrak{A}^k \hat{V} - \mathfrak{A}^k V \right\| \leq \frac{\varepsilon}{3} \tag{C.15}$$

2   Considering the inequality:

$$\left\| \mathfrak{N}^k \hat{V} - \mathfrak{A}^k \hat{V} \right\| \leq \left\| \mathfrak{N}^k - \mathfrak{A}^k \right\| \cdot \left\| \hat{V} \right\| \tag{C.16}$$

3   Obviously, $\left\| \mathfrak{N}^k - \mathfrak{A}^k \right\| \leq M$, and we choose $\left\| \hat{V} \right\| \leq \frac{\varepsilon}{3M}$; it implies:

$$\left\| \mathfrak{N}^k \hat{V} - \mathfrak{A}^k \hat{V} \right\| \leq \frac{\varepsilon}{3} \tag{C.17}$$

4   And we have:

$$\left\| \mathfrak{N}^k V - \mathfrak{A}^k V \right\| \leq \varepsilon \tag{C.18}$$

5