

# Non-linearity matters: a deep learning solution to generalization of hidden brain patterns across population cohorts

Mariam Zabihi<sup>1,2</sup>, Seyed Mostafa Kia<sup>1,2,3</sup>, Thomas Wolfers<sup>1,4,5</sup>, Richard Dinga<sup>1,2</sup>, Alberto Llera<sup>1,2</sup>, Danilo Bzdok<sup>6,7</sup>, Christian F. Beckmann<sup>1,2,8</sup>, Andre Marquand<sup>1,2,9</sup>

<sup>1</sup>Donders Institute for Brain, Cognition and Behavior, Radboud University Nijmegen, Nijmegen, the Netherlands

<sup>2</sup>Department for Cognitive Neuroscience, Radboud University Medical Center Nijmegen, Nijmegen, the Netherlands

<sup>3</sup>Department of Psychiatry, University Medical Center Utrecht, Utrecht, the Netherlands

<sup>4</sup>NORMENT, Division of Mental Health and Addiction, Oslo University Hospital & University of Oslo, Oslo, Norway

<sup>5</sup> Department of Psychology, University of Oslo, Oslo, Norway

<sup>6</sup> Department of Biomedical Engineering, McConnell Brain Imaging Center, Montreal Neurological Institute (MNI), Faculty of Medicine, School of Computer Science, Montreal, Quebec, Canada

<sup>7</sup>Mila - Quebec Artificial Intelligence Institute, Montreal, Quebec, Canada

<sup>8</sup>Centre for Functional MRI of the Brain, University of Oxford, Oxford, United Kingdom

<sup>9</sup>Department of Neuroimaging, Institute of Psychiatry, Psychology, & Neuroscience, King's College London, London, United Kingdom

*Correspondence to:*

Mariam Zabihi

Donders Institute for Brain, Cognition and Behavior

Kapittelweg 29

6525 EN NIJMEGEN

Tel: +3124-3668494

Email: [m.zabihi@donders.ru.nl](mailto:m.zabihi@donders.ru.nl)

## Abstract

The increasing number of neuroimaging scans in recent years has facilitated the use of complex nonlinear approaches to analyzing such data. More specifically, deep learning, which has been previously hindered by the curse of dimensionality is now feasible. However, it remains challenging to use these techniques develop reliable biomarkers and find an optimal representation of data that explains the biological underpinnings of the mental disorders Here, we employed a 3-dimensional autoencoder with an architecture designed from the ground up for task-fMRI data. Our study presented a coherent strategy for optimizing model parameters and architecture and a method for visualizing and interpreting the latent space representation. We trained our model with multi-task fMRI data derived from the Human Connectome Project (HCP) that provides whole-brain coverage across a range of cognitive tasks. Next, in a transfer learning setting, we tested the generalization of our latent space on UK Biobank data as an independent dataset. We showed that the model did not only learn salient features such as age but also high-level behavioral characteristics and that this representation was highly generic and generalizable to an independent dataset. Furthermore, we demonstrated that the projection of latent space back into the original space is meaningful and interpretable. Finally, our results show that with careful implementation, nonlinear features can provide complementary information that accessible to purely linear methods. Our results provide an important step toward learning interpretable and generalizable latent representations that link cognition with underlying brain systems.

## Introduction

The application of machine learning methods and, more specifically, deep-learning methods have shown remarkable returns on many medical imaging problems. [1]. In particular, the applications of deep learning to neuroimaging data are rapidly increasing, mostly the use of supervised learning approaches to solve, for example, classification problems [2]–[9]. However, linear approaches often [10]–[12], but not always [13], produce equal or better performance compared to complex deep neural network models. This is not because the mapping between external covariates and brain correlates is necessarily linear but may instead be due to the lack of data or the choice of input features (e.g., regional averages). While in neuroimaging studies, dealing with a limited number of high-dimensional data had hindered employing such approaches for a time due to the curse of dimensionality [14], the recent increase in the availability of neuroimaging data has provided a great opportunity to move toward employing complex nonlinear methods [12], [13], [15]–[19].

One ultimate challenge in the application of machine learning to neuroimaging is to find reliable biomarkers that explain the biological underpinnings of healthy and disordered mental states. Model interpretability is therefore an important consideration when determining which model to use.[1], [20]–[23]. Many deep learning studies in neuroimaging use hand-crafted features e.g., regions of interest (ROIs) or image-derived phenotypes (IDPs), which are potentially suboptimal for prediction because (i) hand-crafted features may not accurately capture complex structural or functional brain characteristics e.g. overlapping latent representations encoded in the brain, nor their intricate relationships with behavior and (ii) they do not benefit from the strength of deep neural networks in automatically learning the optimal representation from the data (for example using convolutional filters). Particularly in task fMRI studies, which are designed to study mappings from brain activations to cognition and behavior, there are many challenges in understanding the underlying mechanisms, including the extensive heterogeneity across subjects, finding an optimal representation, and a reliable reference to compare the activations [24]–[30]. Consequently, using hand-crafted features potentially leads to losing crucial information relevant, for example, for understanding inter-subject variability [23], [31]. In these scenarios, learning an optimal representation of high-dimensional neuroimaging data rather than – for example – using pre-defined ROIs may enable us to better understand individual variation and more accurately predict clinical and cognitive variables. This representation, also called a latent representation, allows us to reduce the data dimensionality and extract only the essential features from the data. In other words, a latent representation model maps complex and high-dimensional data into a reduced and low-dimensional space[32]. There are two steps to assess the latent representations here: first, whether the derived latent representation shows a stronger association with clinically relevant covariates compared to data in the original space (e.g., mapping from brain to behavioral scores) and further, whether the latent space can be generalized to new data (new brain scans/new participants) which may have a partially different distribution. In the event that this is

proven applicable, then, the knowledge learned from one (big) dataset can be transferred to modeling smaller datasets in a transfer learning paradigm [33].

Most applications of deep learning in neuroscience focus on learning a latent representation that is optimized for a single supervised learning problem, such as predicting age/ sex (e.g. [13] [10]–[12]). However, this may reduce the generalizability of the learned latent representation to other problems. Therefore, we sought to learn a general-purpose latent space that is not bound to a particular task, and instead aims to learn features from the data that are predictive of many different cognitive scores. There have been a number of efforts to that end, e.g. to generate synthetic neuroimaging data [34]–[37]. However, most of these studies evaluate the data representation on the basis of specific measures like reconstruction error. However, this does not necessarily suggest that the latent space presents relevant features, and what is more important is how accurately such representations can predict behavioral or clinical variables. Although linear data-driven transformations like Principle Component Analysis (PCA) and Independent Component Analysis (ICA) [38]–[42] are widely used for feature representation and dimensionality reduction in neuroimaging, these methods often fail to extract complex nonlinear relationships in data. [43], [44].

In this paper, we propose to explore the value of learning a general purpose nonlinear latent space representation of task-fMRI images using a 3-dimensional autoencoder. Autoencoder neural networks provide a powerful tool in various applications in neuroimaging studies, from image segmentation to abnormality detection and latent representation [4], [15], [16], [45]–[48]. Complementary to these approaches, we are interested in automatically learning contextual features using an autoencoder. Briefly, Autoencoder is a deep neural network architecture that consists of two parts an encoder and a decoder. The encoder projects the inputs to a lower-dimensional latent space using a non-linear transformation. The decoder translates back the latent space to the original space by reconstructing the inputs[49]. In contrast to many previous approaches, this does not require the prior specification of nodes or regions of interest, can learn overlapping representations, can use the full range of spatial patterns in the fMRI signal and takes advantage of the strengths of deep learning, for example by learning convolutional filters that capture low-level features of the images.

More specifically, in a fully data-driven approach, we showed that there is useful information about the data in the nonlinear latent space that is not fully captured by a linear data representation and that such information can be extracted using a hierarchical non-linear autoencoder architecture. Here, we employed an autoencoder with an architecture designed from the ground up for task-fMRI data. We trained our model with multi-task fMRI data derived from the Human Connectome Project (HCP)[24] that provides whole-brain coverage across a range of cognitive tasks. Next, in a transfer learning setting, we tested the generalization of our latent space on a UK Biobank dataset[28] after fine-tuning. Our experimental results show that our nonlinear data representation provides a strong foundation for subsequent in-depth analysis and results in more accurate brain-behavior mappings than a commonly used linear approach.

## Methods

An overview of our approach is shown in Figure 1. Here, we used two different data sets. The first data set is task-based fMRI data from HCP [24] S500 release. The second tfMRI data is from the 2020 UK Biobank imaging release[50].

### HCP data

Imaging data: We used tfMRI contrast data from 468 participants in total (187 males and 281 females, Age= 29.2±3.5) from seven different tasks (emotion processing, gambling, language, relational processing, social cognition, motor, working memory) across 86 contrasts which served as the basis in previous brain-imaging work [51], [52]. This yields a total of N≈40K task-fMRI contrasts. The HCP dataset is well suited for this purpose because the task battery covers a wide range of cognitive domains and the neuronal activations associated with the task provide good coverage of the entire brain [25]. The number of participants may vary from task to task; not all the participants have data in all the tasks. While HCP has a large number of samples, the number of participants is relatively small. Therefore, we split data into 5 subsets in a 5-fold cross-validation scheme. The splits are made at the subject level so that each fold contains all the contrasts for a specified set of subjects. In each fold, about 95 participants (20% of the data) were reserved for the test set (N=8K brain scans) and the rest for the training (N=32K brain scans). For each fold, we trained a separate model.

Non-imaging data: the HCP dataset contains various sets of clinical, behavioral, and cognitive tasks, which we use to assess the quality of the latent representation (i.e. how well it predicts behavior). We grouped the scores into 19 categories e.g. sleep quality contains scores relevant to Pittsburgh Sleep Quality Index (PSQI) scores [53] (see supplementary information for the full list of categories). We only included the measures that their scores are available more than half of participants. Moreover, in line with previous studies [54], [55] the measures that had same value for more than 80% of the participants were excluded from further analysis.

### UK Biobank

Imaging data: we used UK Biobank task-fMRI contrast data from 20781 participants and 5 contrasts, in total N≈104K scans (9,860 males, 10,921 females, Age=63.3 ±7.5). The fMRI data derived from UK Biobank uses the same paradigm as the emotion task from the HCP with only minor modifications (e.g. to accommodate shorter run length) [28], [50]. Since UK Biobank provides a larger number of participants than HCP, we trained separate model for each contrast. We randomly selected N=15585 of participants for the train set and 5196 for the test set. All the contrast-models employ the same dataset configuration (the test and train sets).

Non-imaging data: The UK Biobank study provides an extensive number of clinical, behavioral, lifestyle and cognitive scores, which categorized to eight groups e.g., cognitive phenotypes, lifestyle, and mental health (see supplementary information). The exclusion criteria are the same as HCP measures.

## Preprocessing image data and Model architecture selection:

For both datasets we used the volumetrically preprocessed images in standard reference space provided by the respective consortia [56], [57] (for HCP using the ‘minimally processed’ pipeline [56]). Subsequently the scans images were downsampled from 2mm to 3mm voxel resolution to reduce the computational burden then cropped tightly to the whole brain such that the dimension of the image decreased to 56×64×56. The model was trained on the whole-brain contrast images.

We first conducted a pilot study to determine the optimal autoencoder architecture across a wide range of architectures, optimizers, feature scaling approaches with various numbers of hyper-parameters. To do so, we trained the network using half of the HCP training data and tested the performance of the model on reserved test data containing 30 participants (N≈2580 scans). Note that to ensure accurate estimates of generalizability, the reserved test data was used exclusively during pilot study (i.e. we did not re-use these data in subsequent analyses). The details of the models evaluated can be found in the supplementary Table 3. To preserve the details in the image, we choose 3×3×3 kernels for the convolutional layers. More specifically, here, we tested different image normalization (sample-wise and feature-wise of Min-Max and standardization), the number of latent variates (nodes) in the fully connected layer, the number of filters in each layer in addition to different model's optimizer. The reconstruction error of the model reported using mean squared error (MSE) as performance metric of interest.

## Training the model

The final instance of our autoencoder with the optimal architecture from the pilot study was trained using HCP training data with 200 epochs using early-stopping approach to prevent model from overfitting. Next, we used the trained model and transferred the learned parameters from HCP to UK Biobank. We fine-tuned the model with the same hyper-parameters using UK Biobank data. We tested if the convolutional kernels are learning relevant features through the fine-tuning and demonstrated several random kernels’ weights in the supplementary information. This is important to ensure that the autoencoder avoids learning the trivial solution (i.e. a semi-identity function via down-sampling in the encoding phase and up-sampling in the decoding phase). To compare nonlinear transformation to linear, we trained a PCA model with the same number of components as the number of nodes in the autoencoder's fully-connected layer using HCP data. PCA is a natural choice of baseline reference model because it is equivalent (up to a rotation) to a single fully-connected layer autoencoder with a linear activation function [49]. We applied the same PCA model to UK Biobank to compare the results after transfer learning.

## Latent space representation using UMAP

To further evaluate our model, we visualized the latent space using a Uniform Manifold Approximation and Projection (UMAP) approach [58] with two components. UMAP is a manifold learning technique similar to t-distributed stochastic neighbor embedding (t-SNE)

[59] that preserves the local structure of high dimensional data in a nonlinear space. UMAP is superior to tSNE since it better preserves the global structure of data (in addition to its local structure). Furthermore, it is more stable under perturbation or resampling of the data. Moreover, UMAP is relatively fast which is beneficial when tuning the hyperparameters.

Here, we used UMAP to visualize the learned latent space without any preprocessing steps [58] (see supplementary information for more details).

To assist the interpretation of the latent space, we use a simple method to project back the latent spaces in input (i.e., brain) space. To achieve this, we take advantage of the fact that the UMAP algorithm finds clearly separated clusters for the different fMRI contrasts (see results below). Then, for each contrast, we calculated the center of its cluster (i.e., the centroid of K-means clustering) in 2-dimensional UMAP space. We transformed these centroid points to the latent space (using the inverse UMAP transformation) and used the decoder component of the autoencoder to reconstruct the images corresponding to these cluster centers.

## Behavioral Associations

Furthermore, to assess the biological validity of our latent space, we calculated the linear association between clinical and behavioral measures and the reduced latent space for HCP and UK Biobank data.

We used two different approaches because of the different resampling schemes employed in each dataset. For HCP (trained under cross-validation), we regressed non-imaging scores onto the first two UMAP-components using ordinary least squares (OLS) regression in each contrast individually. We used OLS because: (i) the number of test participants in each fold was relatively small in HCP, and (ii) each split resulted in a different latent space, which would make a straight correlation across splits uninterpretable. In the end, we reported the average association (adjusted-r) across separate regression models estimated for each of the five cross-validation splits. In contrast, a single training/test split was used for UK Biobank, so we calculated a straight Pearson correlation between each UMAP component and non-imaging score individually. We repeated the same procedure for behavioral association analysis on the first two PCA components separately.

## Results

### Model selection: autoencoder architecture and hyper-parameter tuning

The optimal model architecture derived from the pilot study is shown in Figure 2. In more detail, the pilot study (Table 3 in supplementary information) suggested that the relatively larger number of filters in the first layers and last layers of the autoencoder results in lower reconstruction error. Moreover, among different image normalization strategies, robust feature-wise normalization was more effective. The number of nodes in the fully-connected layer was another important hyper-parameter, and our results show that the reconstruction performance improves significantly from 10 nodes to 100 nodes while from 50 to 100 nodes, the improvement is small ((Table 3 in supplementary information)).

Table 1: Mean square for PCA model and Autoencoder on HCP data

	MSE	EXPLAINED VARIANCE
PCA(100)	0.706 $\pm$ 0.009	0.331 $\pm$ 0.002
AE(100)	0.741 $\pm$ 0.012	0.301 $\pm$ 0.007

### Autoencoder outperforms PCA in separating tasks and subtasks

The out-of-sample MSE of PCA, shown in Table 1, is comparable with the autoencoder (but for practical purposes equivalent) performance, however as noted above, the association with behavior is of greater importance for the purposes of assessing the quality of the representation.

The scatterplot of the first two components of PCA and two UMAP components of the autoencoder's latent variables for selected contrasts [25] is shown in Figure 3. **Error! Reference source not found.** (See supplementary information for more tasks). As noted above, the UMAP projections show a good separation between tasks.

### Projection the latent representations to brain images

Figure 4 shows the centroids of contrasts that are back-projected from the UMAP latent space to the original brain space. The patterns of activations for these contrasts show an excellent correspondence with the expected task activations as shown in with previous studies (e.g. [25]). For instance, for language task, our projection of latent space to original image space shows the left lateralization which is accord with previous findings [25]

### Association between latent variables and non-imaging covariates

The radar plots in Figure 5 shows the overall association across contrasts per each task in HCP data. These plots indicate that the autoencoder's latent variables have stronger associations with non-imaging variables in comparison to PCA.

Figure 6 shows the same association for the most important contrasts from social, language, gambling and emotion tasks selected [25] (See supplementary information more extensive results). This shows that the autoencoder consistently produces a stronger association with behavior than the PCA (i.e. produces higher correlation in nearly all contrasts).

Parameter transfer from HCP data to UK Biobank The performance of PCA and fine-tuned model on UK Biobank data (Emotion, face-Shape contrast) are shown in (Table 2) indicating that the autoencoder performed equivalently well to PCA in terms of the reconstruction error and explained variance (or slightly better).

Table 2: Mean square for PCA model and Autoencoder on UK Biobank data

	MSE	EXPLAINED VARIANCE
HCP-PCA(100)	0.529	0.33
Fine-tuned AE(100)	0.458	0.38



Figure 7 shows the Manhattan plot of p-value of univariate correlation between non-imaging measures and brain measures; UMAP components of latent space in addition to mean square error. Overall, the autoencoder detected 19 significant Bonferroni-corrected associations with behavior relative to  $3e$  for PCA, demonstrating that after transfer learning to a second independent dataset, the autoencoder also provides a latent space that links more strongly to behavior. Here, PCA principally detects age while autoencoder links to more high-level features e.g., prospective memory, fluid intelligence, length of working weak or mother's age family history. The Manhattan plots of all of contrasts are available in supplementary information.

## Discussion

In this study, we present an optimal 3-D convolutional autoencoder architecture for non-linear transformation of fMRI data to a lower-dimensional yet more informative latent space. We present a coherent strategy for optimizing parameters and architecture of the model and a method to visualize and interpret the learned latent space representation. We showed that our model learned not only salient features that capture age and other sources of population stratification but also high-level behavioral features and that this representation was highly generic and generalized to the UK Biobank population cohort as an independent dataset.

### Task-fMRI data

Regarding the dataset, HCP task-fMRI data enabled us to estimate a generic latent space representation with diverse cognitive tasks and behavioral and clinical scores [24], [25]. Notably, the HCP data provides good whole brain coverage across all the tasks [25]. During the training, this comprehensive mapping allows the autoencoder to learn the various activation patterns across the brain instead of learning specific task-related effects that may be localized to particular brain regions. To validate the generalizability of latent representation of HCP, we used UKB. Complementarily, UK-biobank contains the Hariri faces-shapes emotion task [60], which is similar with emotion task of HCP in a shorter version. The common contrasts provide a great opportunity for further validation of the model and test the across-cohort generalization of the latent space.

### Latent space

The HCP-derived UMAP representation illustrates that the autoencoder can better differentiate between tasks and contrasts compared to a linear model like PCA. As shown in **Error! Reference source not found.**, contrasts are well-separated in the latent space in comparison to the first two components of PCA. This observation demonstrates that the resulting latent space indicated meaningful features in the data. To the contrary, the PCA representation does not separate the different task contrasts as cleanly, potentially leading to a mixing or averaging of different signals. Linear methods like PCA tend to identify prominent information in the data, e.g. age, site, sex. Indeed, these salient information were the only features that were associated with the latent representation derived from PCA in UK Biobank. These features have a substantial effect on data distribution and are easily identified by linear

models, in contrast complex behavioral associations are often missed [21]. Here, using a nonlinear transformation, we aim to identify more complex features. The fact that our nonlinear mapping is more sensitive to the different tasks compared to PCA indicates that the nonlinear latent space can retrieve a different and more informative dimensions in data.

Autoencoders reduce a complex and multidimensional input space into a non-linear combination of latent factors. In contrast with linear decomposition methods such as PCA or supervised learning methods, which are linked to a specific decision problem, the work presented here, provides an effective tool to distill a non-linear representation of images and a powerful and more accurate pipeline to detect anomalies in brain images [5], [15], [61].

Back-projection of latent representation:

The back-projection of latent space shows that we may be able to translate the learned features to a clinical biomarker. For the majority of the contrasts and particularly language (story-math), social (theory of mind) and relational (relational-baseline), the projection of the center of K-means of latent space to the original scan image space were in line with findings in [25]. It is an increasingly challenging problem to figure out how to adapt machine learning algorithms to the discourses we carry on in neuroimaging[62], [63]. In this context, the meaningful projection of the latent space can be viewed as an example of explainable AI in complex models.

### Non-imaging variables relevance

Overall, the latent variables extracted by autoencoder shows stronger association with non-imaging variables compared with PCA. In the HCP dataset, the autoencoder has shown stronger associations with the non-imaging variables in six of the seven tasks, excepting the relational task where PCA showed a stronger association with behavior. Note that the radar plot is the average correlation of the subtasks with each non-imaging variable within the set of clinical and cognitive measures [51]. The detailed correlation map shown in Figure 6 and **Error! Reference source not found.**, indicates that PCA and autoencoder have a comparable performance for captured brain patterns related to categories like clinical and working memory. Importantly, many interesting measures are not associated with the PCA components but are strongly associated with the autoencoder's latent variables. For instance, many social relationships category measures in the Theory of Mind (TOM) contrast show strong correlations with latent variables from the autoencoder but not in the PCA model (dotted circle in Figure 6).

Another example of relevant-behaviors association with latent variables is the reward-punish contrast; while the overall performance in terms of correlation strength sounds similar in two models, the gambling category measures are showing relatively high correlation with latent variables in the autoencoder but not with PCA components.

These results suggest strongly that the our built autoencoder can be adjusted to identify and predict a different set of non-imaging measures and therefore contains potentially

complementary information to linear models. Note that the HCP autoencoder was trained using all the contrasts- in contrast with UK Biobank which was trained only on single Emotion task- hence, the network may only learn the features relevant to all the contrasts and lose the information on contrasts variability. In spite of that, we still can find relevant non-imaging measures associated with the latent variables, showing that the nonlinear manifold may contain interesting information about the data, which can be simply identified through correlation-based correspondence. The association plots show that the derived latent space on UK Biobank data is associated with otherwise challenging-to-predict measures in the lifestyle and the cognitive-phenotype category alongside sex and age. Whilst we recognize the possibility of many confounding effects in large cohorts like UK Biobank [64], our goal was to show that autoencoders produce a meaningful representation for finding brain-behavior associations. Therefore, we consider that a detailed examination of the potential effects of confounding variables was out of the scope of this paper.

### Network architecture

It should be stressed that here, the architectural hyperparameters of the autoencoder are not optimized to deliver the best results concerning reconstruction error but to provide a meaningful nonlinear mapping from input images to the latent space. Nevertheless, the hyperparameter selection procedure has for simplicity been done on the basis of reconstruction error. Some decisions about the network structure have been made before estimating the model. For example, to preserve the morphology of the images and hence better interpretability, we decided to use a 3-D convolutional network [34]–[36], [45]. In order to control order of latent space, we used dense layer in the bottleneck of the autoencoder [49].

We emphasize that we designed our autoencoder with the specific nature of our high-dimensional neuroimaging data in mind and therefore, a number of constraints were imposed on the model beforehand. For example, the networks evaluated were not particularly deep, also to reduce the memory usage and computational complexity, we took advantage of the weight sharing of convolutional layers. Here, we are in search of low-level features that may be translation invariant, but a more important benefit is that the weight sharing enables the networks to be scaled to whole-brain data [65]. The kernel size was set to be  $3 \times 3 \times 3$  to keep the details of the downsampled image scans. Average pooling layers were positioned right after each convolutional layer to ignore the sharp features, reduce the number of parameters and consequently, minimize the chance of overfitting. We relied on the pilot study to select the rest of the model's parameters, such as the number of filters.

While the new trend in neuroimaging is to use linear transformation instead of nonlinear models [12], [66], [67], here, we show that employing nonlinear mappings can add value. We showed that if the nonlinear model's hyperparameters are selected with careful consideration, then we can uncover particular features in the data that are not often revealed by linear models. As a matter of fact, having a solid methodological pipeline is not specific to complex

nonlinear models, but since the chance of over-engineering and hence overfitting is higher than these methods, it is necessary to carefully develop the architectural model design.

Here, we employed a model with a simple reconstruction loss function while the network's primary goal was finding the relevant features to the non-imaging variables. Accordingly, we under-trained our network to preserve the essential aspects of the data, which could be assumed noise from the loss function perspective. It should be noted that our network is entirely unsupervised, and there is no information e.g., age is fed to the loss, although including additional information in the loss is an interesting future direction, for example in a semi-supervised setting that allows the latent space to partially encode particular features of the data [15]. Another interesting future direction is to train an autoencoder to predict different data (e.g., a follow-up timepoint in longitudinal studies). This would serve to sensitize the latent space to changes relevant to ageing or pathology.

The increased number of neuroimaging scans provides a unique opportunity to transcend the linear mapping, but it is also necessary to acknowledge some limitations. The traditional image processing techniques are not completely applicable here. For example, while data augmentation using image mirroring, flipping, skewing, or segmenting is a straightforward approach to increase the number of samples and has been applied before in neuroimaging applications [13], we did not consider it to be appropriate here because such augmentation strategies do not faithfully preserve invariances known to occur in the brain, for example the lateralization of brain functions e.g. the association of left lateralization in language processing [68]. Another limitation is computational complexity. Training an autoencoder on large neuroimaging data is computationally more demanding comparing with similar linear models.

## Conclusion

This study aimed to explore nonlinear manifold in brain imaging using 3-dimensional autoencoders. We presented a coherent strategy for optimizing parameters and architecture of the model and a method to visualize and interpret the learned latent space representation. We showed that our model learned not only salient features like age but also high-level behavioral features and that this representation was highly generic and generalized to an independent dataset. Finally, our results show that with a careful implementation, nonlinear features can provide complementary information that accessible to purely linear methods.

## References

- [1] G. Litjens *et al.*, “A survey on deep learning in medical image analysis,” *Med. Image Anal.*, vol. 42, no. December 2012, pp. 60–88, 2017.
- [2] D. Bzdok and A. Meyer-Lindenberg, “Machine Learning for Precision Psychiatry: Opportunities and Challenges,” *Biological Psychiatry: Cognitive Neuroscience and Neuroimaging*, vol. 3, no. 3. Elsevier Inc, pp. 223–230, 01-Mar-2018.
- [3] D. Durstewitz, G. Koppe, and A. Meyer-Lindenberg, “Deep neural networks in psychiatry,” *Mol. Psychiatry*, p. 1, Feb. 2019.
- [4] H. Il Suk, S. W. Lee, and D. Shen, “Deep ensemble learning of sparse regression models for brain disease diagnosis,” *Med. Image Anal.*, vol. 37, pp. 101–113, Apr. 2017.
- [5] S. Vieira, W. H. L. Pinaya, and A. Mechelli, “Using deep learning to investigate the neuroimaging correlates of psychiatric and neurological disorders: Methods and applications,” *Neurosci. Biobehav. Rev.*, vol. 74, pp. 58–75, Mar. 2017.
- [6] T. Wolfers, J. K. Buitelaar, C. F. Beckmann, B. Franke, and A. F. Marquand, “From estimating activation locality to predicting disorder: A review of pattern recognition for neuroimaging-based psychiatric diagnostics,” *Neurosci. Biobehav. Rev.*, vol. 57, pp. 328–349, 2015.
- [7] A. Mensch, J. Mairal, D. Bzdok, B. Thirion, and G. Varoquaux, “Learning Neural Representations of Human Cognition across Many fMRI Studies,” no. Nips, 2017.
- [8] M. R. Arbabshirani, S. Plis, J. Sui, and V. D. Calhoun, “Single subject prediction of brain disorders in neuroimaging: Promises and pitfalls,” *Neuroimage*, vol. 145, pp. 137–165, 2017.
- [9] T. Wolfers *et al.*, “From pattern classification to stratification: towards conceptualizing the heterogeneity of Autism Spectrum Disorder,” *Neuroscience and Biobehavioral Reviews*, vol. 104. Elsevier Ltd, pp. 240–254, 01-Sep-2019.
- [10] C. Davatzikos, “Machine learning in neuroimaging: Progress and challenges,” *Neuroimage*, vol. 197, p. 652, 2019.
- [11] T. He *et al.*, “Deep neural networks and kernel regression achieve comparable accuracies for functional connectivity prediction of behavior and demographics,” *Neuroimage*, vol. 206, Feb. 2020.
- [12] M. A. Schulz *et al.*, “Different scaling of linear models and deep learning in UKBiobank brain images versus machine-learning datasets,” *Nat. Commun.*, vol. 11, no. 1, 2020.
- [13] H. Peng, W. Gong, C. F. Beckmann, A. Vedaldi, and S. M. Smith, “Accurate brain age prediction with lightweight deep neural networks,” *Med. Image Anal.*, vol. 68, p. 101871, Feb. 2021.
- [14] R. E. Bellman, *Adaptive control processes: a guided tour*, vol. 2045. Princeton university press, 2015.
- [15] W. H. L. Pinaya, A. Mechelli, and J. R. Sato, “Using deep autoencoders to identify abnormal brain structural patterns in neuropsychiatric disorders: A large-scale multi-

- sample study," *Hum. Brain Mapp.*, vol. 40, no. 3, pp. 944–954, Feb. 2019.
- [16] W. H. L. Pinaya *et al.*, "Normative modelling using deep autoencoders: a multi-cohort study on mild cognitive impairment and Alzheimer's disease," *bioRxiv*, 2020.
  - [17] J. H. Cole *et al.*, "Predicting brain age with deep learning from raw imaging data results in a reliable and heritable biomarker," *Neuroimage*, vol. 163, pp. 115–124, 2017.
  - [18] N. K. Dinsdale *et al.*, "Learning patterns of the ageing brain in MRI using deep convolutional networks," *Neuroimage*, vol. 224, p. 117401, Jan. 2021.
  - [19] H. Kiesow *et al.*, "Hidden population modes in social brain morphology: Its parts are more than its sum," *bioRxiv*, p. 2020.08.07.241497, Aug. 2020.
  - [20] R. A. Poldrack *et al.*, "Scanning the horizon: Towards transparent and reproducible neuroimaging research," *Nat. Rev. Neurosci.*, vol. 18, no. 2, pp. 115–126, 2017.
  - [21] B. Mwangi, T. S. Tian, and J. C. Soares, "A review of feature reduction techniques in Neuroimaging," *Neuroinformatics*, vol. 12, no. 2. Humana Press Inc., pp. 229–244, 2014.
  - [22] S. Korolev, A. Safiullin, M. Belyaev, and Y. Dodonova, "Residual and plain convolutional neural networks for 3D brain MRI classification," *Proc. - Int. Symp. Biomed. Imaging*, pp. 835–838, 2017.
  - [23] W. Gong, C. F. Beckmann, and S. M. Smith, "Phenotype Discovery from Population Brain Imaging," *bioRxiv*, 2020.
  - [24] D. C. Van Essen, S. M. Smith, D. M. Barch, T. E. J. Behrens, E. Yacoub, and K. Ugurbil, "The WU-Minn Human Connectome Project: An overview," *Neuroimage*, vol. 80, pp. 62–79, 2013.
  - [25] D. M. Barch *et al.*, "Function in the human connectome: Task-fMRI and individual differences in behavior," *Neuroimage*, vol. 80, pp. 169–189, 2013.
  - [26] S. M. Smith *et al.*, "A positive-negative mode of population covariation links brain connectivity, demographics and behavior," *Nature Neuroscience*, vol. 18, no. 11. Nature Publishing Group, pp. 1565–1567, 01-Nov-2015.
  - [27] E. S. Finn *et al.*, "Functional connectome fingerprinting: identifying individuals using patterns of brain connectivity," *Nat. Neurosci.*, vol. 18, no. October, pp. 1–11, 2015.
  - [28] K. L. Miller *et al.*, "Multimodal population brain imaging in the UK Biobank prospective epidemiological study," *Nat. Neurosci.*, vol. 19, no. 11, pp. 1523–1536, 2016.
  - [29] L. Gupta *et al.*, "Spatial heterogeneity analysis of brain activation in fMRI," *NeuroImage Clin.*, vol. 5, pp. 266–276, 2014.
  - [30] G. C. Burgess, J. R. Gray, A. R. A. Conway, and T. S. Braver, "Neural mechanisms of interference control underlie the relationship between fluid intelligence and working memory span," *J. Exp. Psychol. Gen.*, 2011.
  - [31] D. Bzdok, "Classical statistics and statistical learning in imaging neuroscience," *Frontiers in Neuroscience*, vol. 11, no. OCT. Frontiers Media S.A., p. 543, 06-Oct-2017.

- [32] H. Il Suk, S. W. Lee, and D. Shen, "Latent feature representation with stacked auto-encoder for AD/MCI diagnosis," *Brain Struct. Funct.*, vol. 220, no. 2, pp. 841–859, 2015.
- [33] Y. Bengio, "Deep learning of representations for unsupervised and transfer learning," in *Proceedings of ICML workshop on unsupervised and transfer learning*, 2012, pp. 17–36.
- [34] P.-D. Tudosiu *et al.*, "Neuromorphologically-preserving Volumetric data encoding using VQ-VAE," *arXiv*, pp. 1–13, Feb. 2020.
- [35] G. Kwon, C. Han, and D. Kim, "Generation of 3D Brain MRI Using Auto-Encoding Generative Adversarial Networks," *Lect. Notes Comput. Sci. (including Subser. Lect. Notes Artif. Intell. Lect. Notes Bioinformatics)*, vol. 11766 LNCS, pp. 118–126, Aug. 2019.
- [36] H. Choi, H. Kang, and D. S. Lee, "Predicting aging of brain metabolic topography using variational autoencoder," *Front. Aging Neurosci.*, vol. 10, no. JUL, p. 212, Jul. 2018.
- [37] H. Huang *et al.*, "Modeling Task fMRI Data Via Deep Convolutional Autoencoder," *IEEE Trans. Med. Imaging*, vol. 37, no. 7, pp. 1551–1561, Jul. 2018.
- [38] F. Bunea, Y. She, H. Ombao, A. Gongvatana, K. Devlin, and R. Cohen, "Penalized least squares regression methods and applications to neuroimaging," *Neuroimage*, vol. 55, no. 4, pp. 1519–1527, Apr. 2011.
- [39] G. Sidhu, N. Asgarian, R. Greiner, and M. R. G. Brown, "Kernel principal component analysis for dimensionality reduction in fMRI-based diagnosis of ADHD," *Front. Syst. Neurosci.*, vol. 6, no. October, pp. 1–17, Oct. 2012.
- [40] G. Varoquaux, P. R. Raamana, D. A. Engemann, A. Hoyos-Idrobo, Y. Schwartz, and B. Thirion, "Assessing and tuning brain decoders: Cross-validation, caveats, and guidelines," *Neuroimage*, vol. 145, pp. 166–179, Jan. 2017.
- [41] V. D. Calhoun, J. Liu, and T. Adali, "A review of group ICA for fMRI data and ICA for joint inference of imaging, genetic, and ERP data," *Neuroimage*, vol. 45, no. 1 Suppl, 2009.
- [42] B. Thirion and O. Faugeras, "Dynamical components analysis of fMRI data through kernel PCA," *Neuroimage*, vol. 20, no. 1, pp. 34–49, Sep. 2003.
- [43] D. Bzdok and B. T. T. Yeo, "Inference in the age of big data: Future perspectives on neuroscience," *NeuroImage*, vol. 155. Academic Press Inc., pp. 549–564, 15-Jul-2017.
- [44] S. M. Smith and T. E. Nichols, "Statistical Challenges in 'Big Data' Human Neuroimaging," *Neuron*, vol. 97, no. 2. Cell Press, pp. 263–268, 17-Jan-2018.
- [45] A. Payan and G. Montana, "Predicting Alzheimer 's disease : a neuroimaging study with 3D convolutional neural networks," *arXiv Prepr. arXiv1502.02506*, pp. 1–9, 2015.
- [46] J. E. Savage *et al.*, "Genome-wide association meta-analysis in 269,867 individuals identifies new genetic and functional links to intelligence," *Nat. Genet.*, vol. 50, no. 7, pp. 912–919, Jul. 2018.
- [47] H. Huang *et al.*, "Modeling Task fMRI Data Via Deep Convolutional Autoencoder," *IEEE*



- Trans. Med. Imaging*, vol. 37, no. 7, pp. 1551–1561, Jul. 2018.
- [48] A. Myronenko, “3D MRI brain tumor segmentation using autoencoder regularization,” in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2019, vol. 11384 LNCS, pp. 311–320.
  - [49] I. Goodfellow, Y. Bengio, A. Courville, and Y. Bengio, *Deep learning*, vol. 1, no. 2. MIT press Cambridge, 2016.
  - [50] T. J. Littlejohns *et al.*, “The UK Biobank imaging enhancement of 100,000 participants: rationale, data collection, management and future directions,” *Nat. Commun.*, vol. 11, no. 1, p. 2624, 2020.
  - [51] D. Bzdok, G. Varoquaux, O. Grisel, M. Eickenberg, C. Poupon, and B. Thirion, “Formal Models of the Network Co-occurrence Underlying Mental Operations,” *PLOS Comput. Biol.*, vol. 12, no. 6, p. e1004994, Jun. 2016.
  - [52] D. Bzdok, M. Eickenberg, O. Grisel, B. Thirion, G. Varoquaux Semi, and G. Varoquaux, “Semi-Supervised Factored Logistic Regression for High-Dimensional Neuroimaging Data,” 2015.
  - [53] D. J. Buysse, C. F. Reynolds, T. H. Monk, S. R. Berman, and D. J. Kupfer, “The Pittsburgh sleep quality index: A new instrument for psychiatric practice and research,” *Psychiatry Res.*, vol. 28, no. 2, pp. 193–213, 1989.
  - [54] S. M. Smith *et al.*, “A positive-negative mode of population covariation links brain connectivity, demographics and behavior,” *Nat. Neurosci.*, vol. 18, no. 11, pp. 1565–1567, 2015.
  - [55] A. F. Marquand, K. V. Haak, and C. F. Beckmann, “Functional corticostriatal connection topographies predict goal-directed behaviour in humans,” *Nat. Hum. Behav.*, vol. 1, no. 8, p. 146, Jul. 2017.
  - [56] M. F. Glasser *et al.*, “The minimal preprocessing pipelines for the Human Connectome Project,” *Neuroimage*, vol. 80, pp. 105–124, Oct. 2013.
  - [57] F. Alfaro-Almagro *et al.*, “Image processing and Quality Control for the first 10,000 brain imaging datasets from UK Biobank,” *Neuroimage*, vol. 166, pp. 400–424, 2018.
  - [58] L. McInnes, J. Healy, and J. Melville, “UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction,” *arXiv*, Feb. 2018.
  - [59] L. der Maaten and G. Hinton, “Visualizing data using t-SNE,” *J. Mach. Learn. Res.*, vol. 9, no. 11, 2008.
  - [60] A. R. Hariri, A. Tessitore, V. S. Mattay, F. Fera, and D. R. Weinberger, “The amygdala response to emotional stimuli: A comparison of faces and scenes,” *Neuroimage*, vol. 17, no. 1, pp. 317–323, 2002.
  - [61] P. Perera and V. M. Patel, “Learning deep features for one-class classification,” *IEEE Trans. Image Process.*, vol. 28, no. 11, pp. 5450–5463, 2019.
  - [62] B. Heinrichs and S. B. Eickhoff, “Your evidence? Machine learning algorithms for medical diagnosis and prediction,” *Hum. Brain Mapp.*, vol. 41, no. 6, pp. 1435–1444,



Apr. 2020.

- [63] D. Bzdok and J. P. A. Ioannidis, “Exploration, inference, and prediction in neuroscience and biomedicine,” *Trends Neurosci.*, vol. 42, no. 4, pp. 251–262, 2019.
- [64] F. Alfaro-Almagro *et al.*, “Confound modelling in UK Biobank brain imaging,” *Neuroimage*, vol. 224, p. 117002, 2021.
- [65] Y. Lecun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, vol. 521, no. 7553. Nature Publishing Group, pp. 436–444, 27-May-2015.
- [66] A. D’amour *et al.*, “Underspecification in Machine Learning Underspecification Presents Challenges for Credibility in Modern Machine Learning,” 2020.
- [67] A. Saxe, S. Nelli, and C. Summerfield, “If deep learning is the answer, what is the question?,” *Nat. Rev. Neurosci.*, 2020.
- [68] J. A. Frost *et al.*, “Language processing is strongly left lateralized in both sexes. Evidence from functional MRI,” *Brain*, vol. 122, no. 2, pp. 199–208, 1999.

Figures:

Figure 1: Method's overview

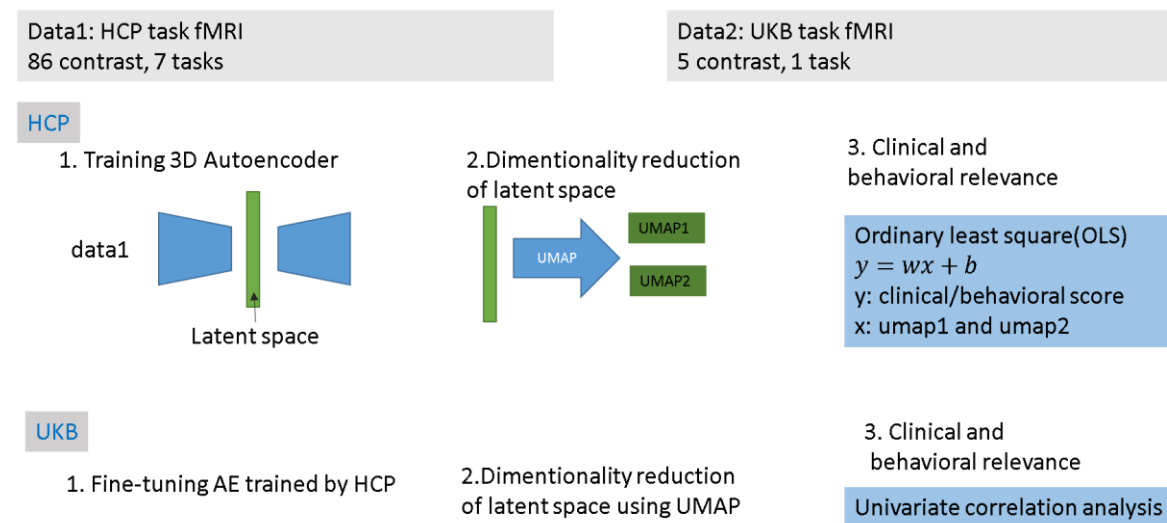


Figure 2: Autoencoder architecture

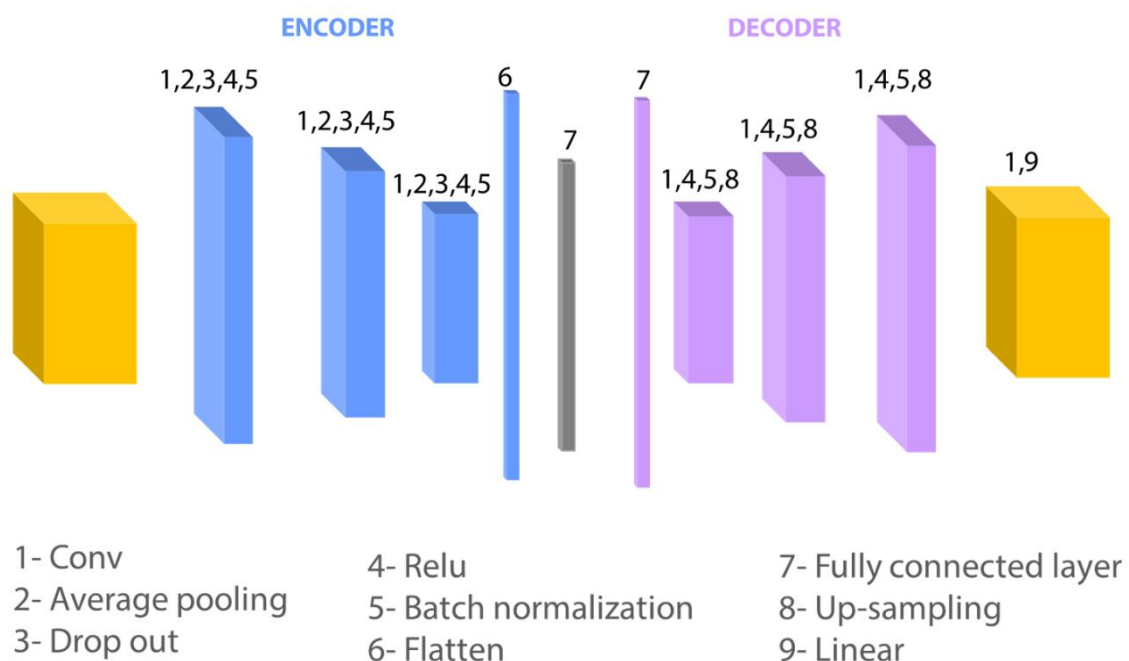


Figure 3: UMAP of the latent space versus two PCA components for selected contrasts according to Barch 2013

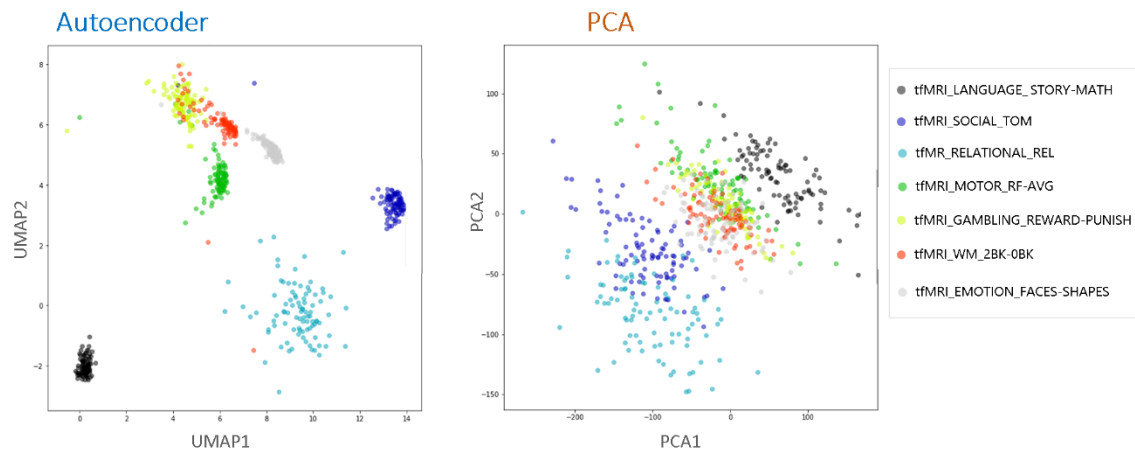


Figure 4: The projection of centroid of UMAP in the latent space to the input brain space. The centers of UMAP of latent space were calculated using K-means clustering across the test data. The centroids corresponding to each contrast were passed to encoder of autoencoder to map to input original space.

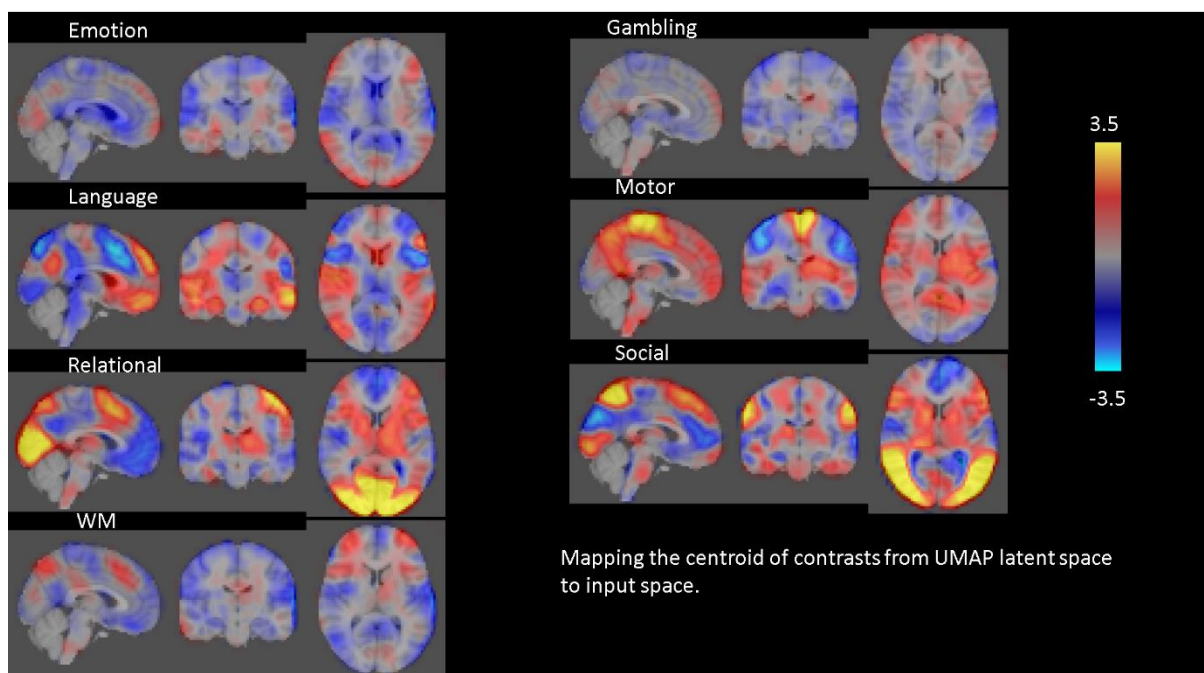


Figure 5: The summary of association of latent variables and non-imaging measures across different behavioral categories in HCP data. Note that this is the average of association (OLS adjusted-r) over five separated models. Autoencoder indicates stronger associations in most of targets compare to PCA.

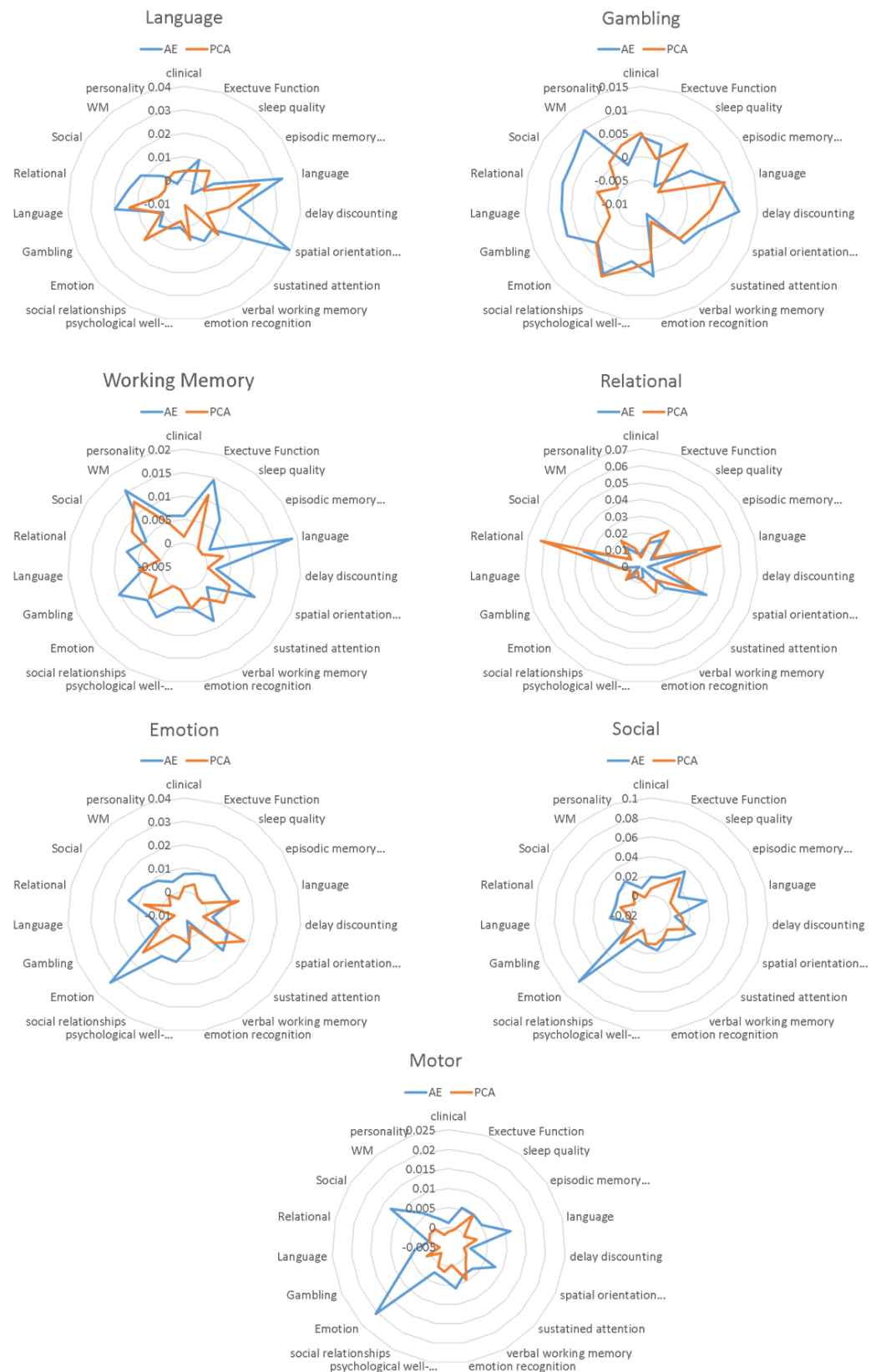


Figure 6: The association of UMAP of latent variables and non-brain imaging variables in HCP contrasts. Showing contrasts were selected according to [25]. The association (adjusted-r score of OLS) of each contrast were averaged across five-fold cross-validation. The dotted circles show the higher associations for some behavioral measures in the relevant task in the autoencoder model compared to PCA.

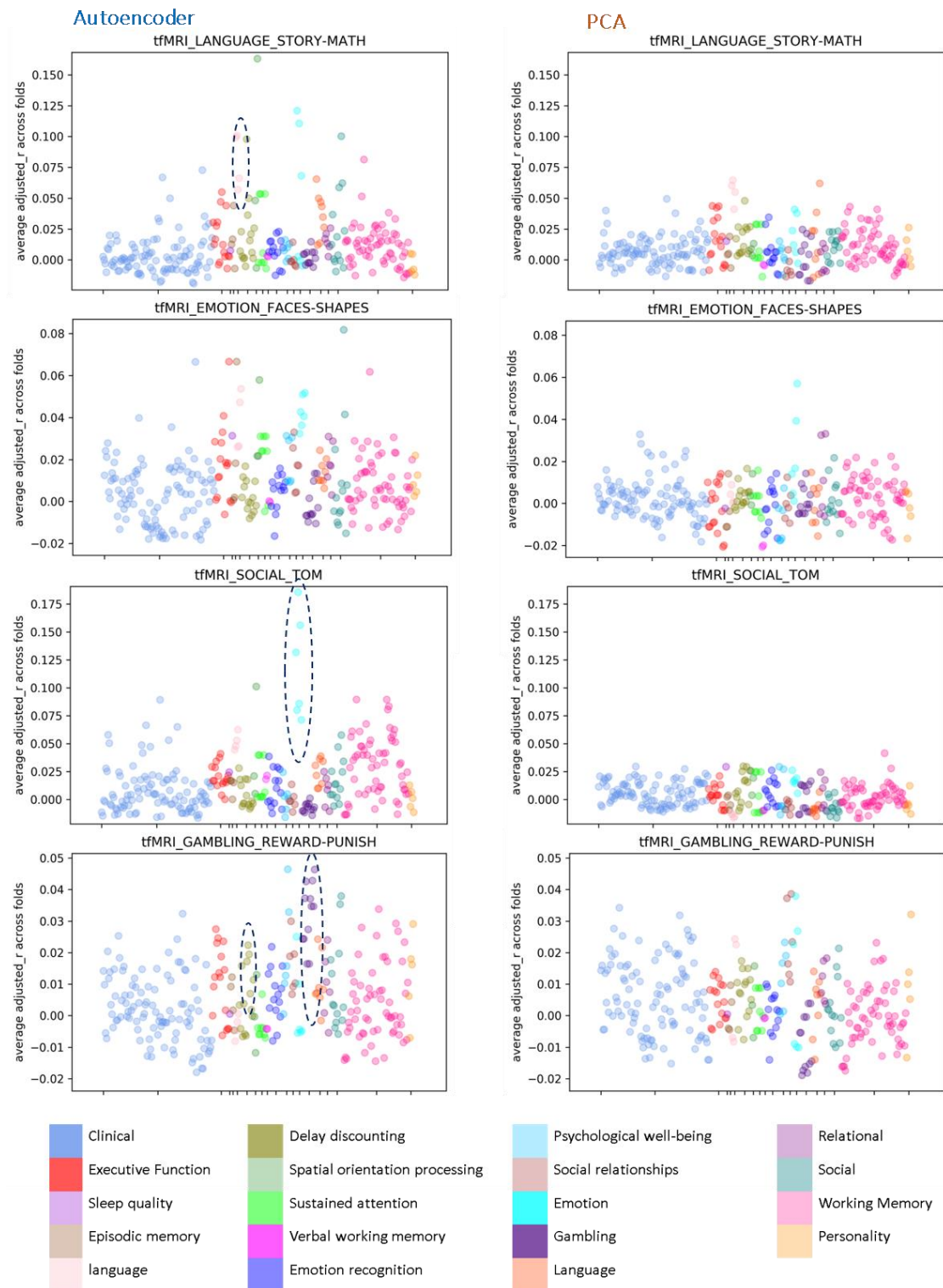


Figure 7: Manhattan plot of  $p$ -value of univariate correlation of non-imaging measures with UMAP of latent variables (left) and first two components of PCA (right) for faces-Shapes subtask. The black line is Bonferroni-corrected  $p$ -value threshold.

