

DemoVAE

A Demographic-Conditioned Variational Autoencoder Reveals Most fMRI-Based Prediction Depends on Demographic Confounds

Anton Orlichenko¹, Gang Qu¹, Ziyu Zhou², Anqi Liu³, Hong-Wen Deng³, Zhengming Ding², Vince Calhoun⁴, Yu-Ping Wang^{1,2}

¹Department of Biomedical Engineering, Tulane University, New Orleans, LA

²Department of Computer Science, Tulane University, New Orleans, LA

³Center for Biomedical Informatics and Genomics, Tulane University, New Orleans, LA

⁴TReNDS, Georgia State University, Atlanta, GA

May 1, 2024



Background: fMRI

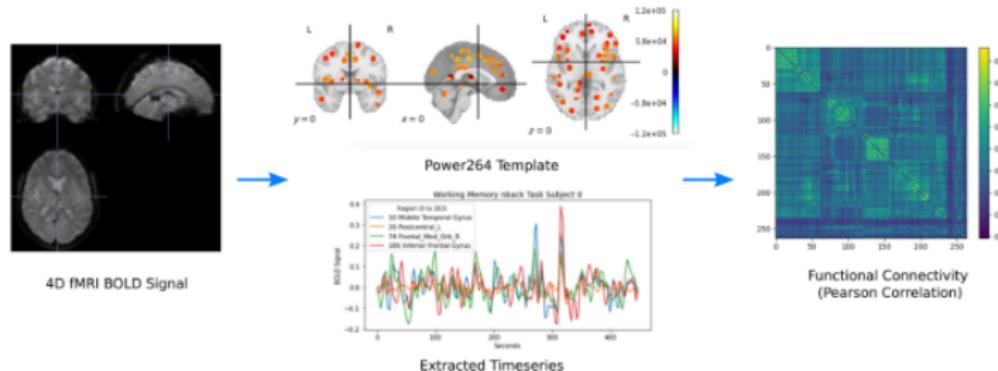


Fig. 1 Preprocessing pipeline for converting 4D fMRI volumes into FC using the Power264 atlas.³⁴
Reproduced with permission from Orlichenko et al.³⁵

- fMRI is a non-invasive measure of neural activity via the blood oxygen level-dependent (BOLD) signal
- Functional connectivity (FC) is the temporal Pearson correlation of BOLD signal between different regions in the brain
- The brain is usually parcellated according to an atlas



Background: Prediction of Clinical Features Using fMRI

TABLE V
RESULTS OF PNC EXPERIMENTS.

Model	Age (RMSE, years)		Sex (Accuracy)		Intelligence (RMSE, WRAT score)	
	N=30	N=496	N=30	N=496	N=30	N=496
Null	3.3	0.54	15.7			
M-GCN	4.47	2.37	0.51	0.75	23.27	15.59
MLP	4.52	2.43	0.53	0.8	21.17	15.64
GCN	3.89	2.16	0.49	0.8	16.29	14.38
Linear	4.36	2.34	0.54	0.77	19.8	14.97
LatSim	2.86	2.05	0.55	0.82	15.59	14.26
p-value	2.2e-6	5e-3	0.32	0.11	0.02	0.30

Table 4: Schizophrenia classification results on the CNP dataset.

Model	Accuracy	AUC	BAC
TFF	88.2	90.0	87.9
TFF_vanilla	58.8	52.9	50.0
ST-GCN	82.3	87.1	82.8
Deep-fMRI	76.5	84.3	77.9

Table 5: Ablation study, see text for variants (i)–(v).

	Age pred.			Gender pred.		
	L1	L2	NMSE	Acc.	BAC	AUC
(i) no intensity loss	2.96	11.71	0.16	93.77	92.95	96.06
(ii) no L1 loss	3.12	13.07	0.18	89.66	87.43	91.25
(iii) no perceptual loss	3.02	12.11	0.17	93.47	93.02	96.86
(iv) no 2-norm	3.09	12.65	0.17	93.07	92.96	93.37
(v) one-step pre-training	3.21	14.33	0.19	89.97	91.02	90.38
Full method	2.73	10.93	0.14	94.09	93.92	98.77

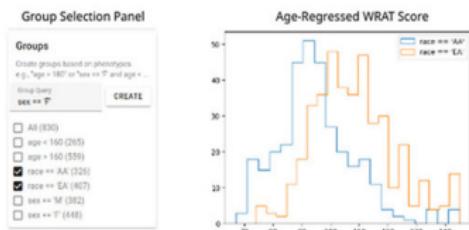
- FC has been used as input to various models in order to predict age, sex, race, general fluid intelligence, and disease status¹
- It is also possible to use 4D fMRI data directly to perform prediction with CNN or transformer architecture²

¹Orlichenko et al. 2022 <https://doi.org/10.1109/TBME.2022.3232964>

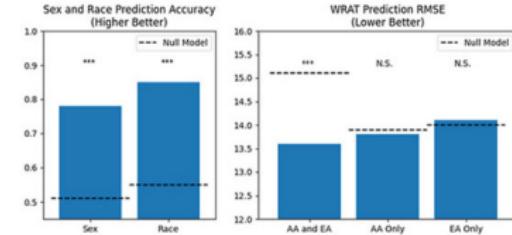
²Malkiel et al. 2022 <https://doi.org/10.48550/arXiv.2112.05761>

Background: Previous Work On Demographic Confounds

C. Finds Ethnic Bias in WRAT Score



D. FC Prediction of WRAT Actually Predicts Race



- Others³ as well as ourselves⁴ have shown that achievement score or general fluid intelligence prediction may be confounded by different score distributions between races

³ Li et al. 2022 <https://doi.org/10.1126/sciadv.abj1812>

⁴ Orlichenko et al. 2023 <https://doi.org/10.1016/j.ynirp.2023.100191>

Motivation

Question of Demographic Confounds

- Previous work has shown demographics confound general intelligence prediction
- To what extent do demographics influence other fMRI-based prediction?

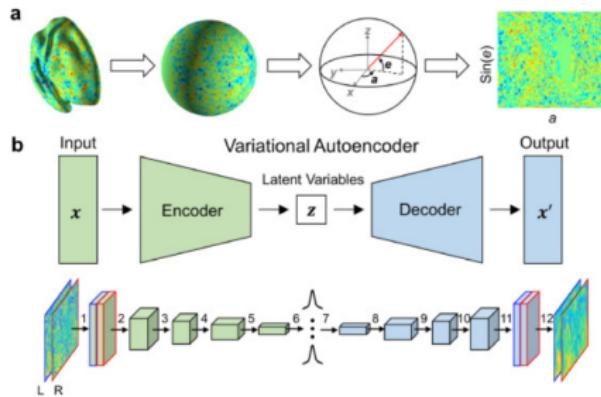
Synthetic Data for Broader Access

- Access to large fMRI datasets often requires application by qualified researchers and IRB approval
- Could we create a method by which a model could "learn" a dataset and be distributed without compromising patient privacy?

Data Harmonization

- Can we remove site-specific effects from fMRI data?

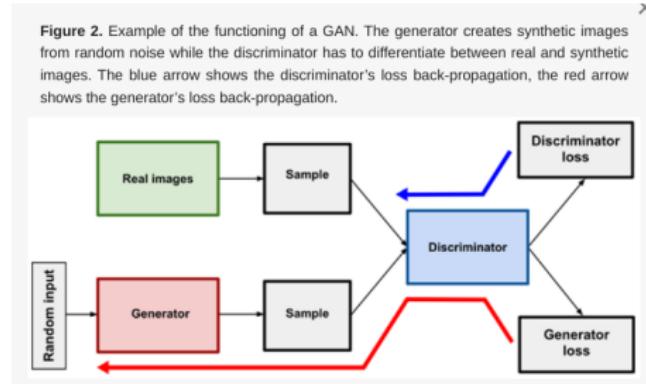
Background: VAEs and fMRI



- Previous works have used variational autoencoders (VAEs) to model fMRI BOLD activity⁵
- A VAE learns to convert input samples to a latent space with a fixed distribution, allowing for later sampling from the latent space

⁵ Kim et al. 2021 <https://doi.org/10.1016/j.neuroimage.2021.118423>

Background: Generative Adversarial Networks



- Other studies have used Generative Adversarial Networks (GANs)
- GANs use a generator module to create samples from noise and a discriminator module to teach the generator to create realistic samples



Background: Generative Models in Computer Vision



- Well-known apps and models such as DALLE-2, Midjourney, and Stable Diffusion use (or have used) VAEs, GANs, and diffusion-based models to create synthetic images which appear almost entirely realistic⁶

⁶ Image credit: NVidia

<https://venturebeat.com/ai/gan-generative-adversarial-network-explainer-ai-machine-learning/>

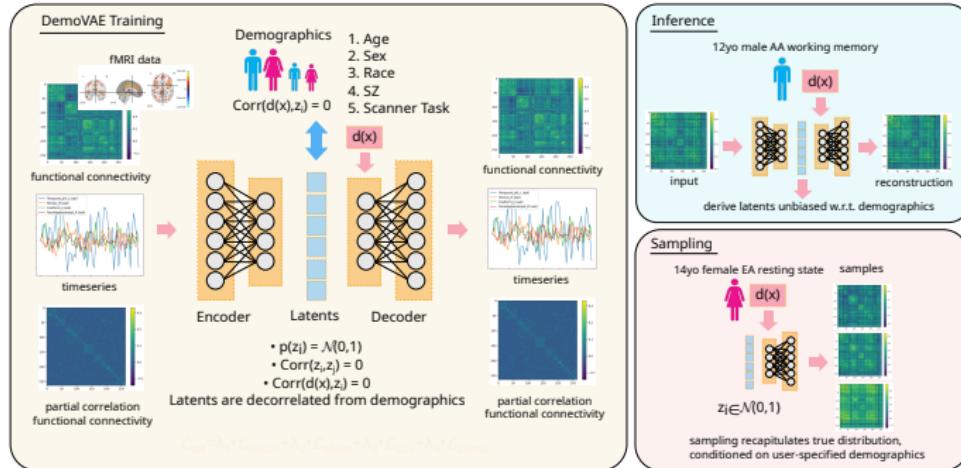
Background: Topics Explored

- The topics explored by previous GAN studies include:⁷
 - brain decoding (4 studies)
 - disease progression modeling (2 studies)
 - image reconstruction (9 studies)
 - image segmentation (4 studies)
 - image synthesis (7 studies)
 - image to image translation (18 studies)
- We believe the conditioning of generative models on demographics in order to explore the effects of demographic confounds and create unconfounded latents is a potentially useful work



⁷ Laino et al. 2022 <https://doi.org/10.3390/jimaging8040083>

Our Solution



- We create a VAE which encodes fMRI FC data and is conditioned on patient demographics
- This requires de-correlating the latent state z from demographic features y

VAE Theory: KL Divergence

- The VAE is a generative model which encodes an input x into a latent space with a known Gaussian distribution $p_\theta(z) = \mathcal{N}(0, 1)$
- This is done by minimizing the Kullback-Leibler (KL) divergence between the true conditional distribution $p_\theta(z|x)$ and $q_\phi(z|x)$:

$$\begin{aligned} D_{KL}(q_\phi(z|x) \| p_\theta(z|x)) &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\ln \frac{q_\phi(z|x)}{p_\theta(z|x)} \right] \\ &= \ln p_\theta(x) + \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\ln \frac{q_\phi(z|x)}{p_\theta(x, z)} \right] \end{aligned} \tag{1}$$



VAE Theory: Evidence Lower Bound

- The Evidence Lower Bound (ELBO) is then defined
- The ELBO has two parts: a reconstruction fidelity part to be maximized and a KL divergence part to be minimized⁸

$$\begin{aligned} L_{\theta, \phi} &= \mathbb{E}_{z \sim q_{\phi}(\cdot|x)} \left[\ln \frac{p_{\theta}(x, z)}{q_{\phi}(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_{\phi}(\cdot|x)} [\ln p_{\theta}(x|z) - D_{KL}(q_{\phi}(z|x) \| p_{\theta}(z))] \end{aligned} \tag{2}$$



⁸ Kingma et al. 2014 <https://doi.org/10.48550/arXiv.1312.6114>

Scalar or Vector Latent Features

- In the case of scalar latent features being distributed according to $p_\theta(z) = \mathcal{N}(0, 1)$, maximizing the ELBO is equivalent to minimizing the following loss:

$$\mathcal{L}_{\theta, \phi} = \|x - D_\theta(z)\|_2^2 + N\sigma_z^2 + \|\mu_z\|_2^2 - N\ln\sigma_z^2, \quad (3)$$

- In the case of multidimensional latent features $p_\theta(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$, the KL divergence part takes a more complicated form:⁹

$$\begin{aligned} D_{KL}(\mathcal{N}(\mu_z, \Sigma_z) || \mathcal{N}(\mathbf{0}, \mathbf{I})) &= \\ \frac{1}{2} \left[\text{tr}(\Sigma_z) + \mu_z^\top \mu_z + \log(\det(\Sigma_z)) \right] \end{aligned} \quad (4)$$

- This is problematic because of the log determinant term



⁹ Murphy, Kevin P. 2023 <http://probml.github.io/book2>

Modification of Training Objective

- We modify the regular VAE training objective with new terms for reconstruction and KL divergence
- We also add several terms important for conditioning the VAE output on demographics
 - ① Incorporate demographic information into reconstruction
 - ② Extend to multidimensional latent space
 - ③ Decorrelate latent features from demographics
 - ④ Add classifier guidance



1. Incorporate Demographic Information

- We modify the decoder $\hat{\mathbf{x}} = D_\theta(\mathbf{z})$ used for sample reconstruction to be conditioned on user-specified demographics ($\hat{\mathbf{x}} = D_\theta(\mathbf{z}, \mathbf{y})$)

$$\mathcal{L}_{\text{Recon}} = \frac{1}{N} \sum_{i=1}^N \|\mathbf{x}_i - D_\theta(\mathbf{z}_i, \mathbf{y}_i)\|_2^2, \quad (5)$$



2. Extension to Multidimensional Latent Space

- We replace the log determinant-containing KL divergence term with terms which match the multidimensional mean (zero) and covariance (diagonal \mathbf{I} matrix) of a multivariate $\mathcal{N}(\mathbf{0}, \mathbf{I})$ distribution

$$\begin{aligned}\mathcal{L}_{\text{Cov}} &= \frac{1}{N} \|\mathbf{Z}\mathbf{Z}^\top - N\mathbf{I}\|_F^2 \\ \mathcal{L}_{\text{Mean}} &= \frac{1}{NN_z} \sum_{i=1}^{N_z} \|\mu_{\mathbf{z}_i}\|_2^2,\end{aligned}\tag{6}$$

3. Decorrelate Latent Features from Demographics

- All demographic information in final output $\hat{\mathbf{x}} = D_\theta(\mathbf{z}, \mathbf{y})$ should come from user-specified demographics \mathbf{y} and not the latent features \mathbf{z}
- In order to achieve this, we force the model to create latent features which are decorrelated from demographics during training

$$\mathcal{L}_{\text{Demo}} = \frac{1}{N_z N_y} \sum_{j=1}^{N_z} \sum_{k=1}^{N_y} \|\rho_{\mathbf{z}_j, \mathbf{y}_k}\|_2^2, \quad (7)$$



4. Classifier Guidance

- Finally, we force synthetic FC data generated by DemoVAE conditioned upon user-specified demographics to match what is expected by models trained on real data to predict those demographic features

$$\mathcal{L}_{\text{Guide}} = \frac{1}{N N_y} \sum_{i=1}^{N_y} \begin{cases} \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2^2, & \mathbf{y}_i \text{ continuous} \\ -\sum_c \mathbf{y}_{i,c} \log(p_{i,c}), & \mathbf{y}_i \text{ categorical} \end{cases} \quad (8)$$

Total Training Objective

- The total loss function is the following:

$$\mathcal{L} = \mathcal{L}_{\text{Recon}} + \lambda_1 \mathcal{L}_{\text{Cov}} + \lambda_2 \mathcal{L}_{\text{Mean}} + \lambda_3 \mathcal{L}_{\text{Demo}} + \lambda_4 \mathcal{L}_{\text{Guide}}, \quad (9)$$

- where λ_{1-4} are hyperparameters that are chosen using random grid search.

Datasets

- We train and validate on two large, widely used datasets available through application to NIH
- Philadelphia Neurodevelopmental Cohort (PNC)
 - Dataset of children and young adults 8-22 years old
 - Three fMRI scanner tasks: resting state, working memory (nback), and emotion identification (emoid)
 - 169 questionnaire and computerized battery fields
 - fMRI data for 1,529 subjects and SNP data for 9,000+ subjects
 - 1,154 subjects have all three scanner tasks and SNP data
- Bipolar and Schizophrenia Network for Intermediate Phenotypes (B-SNIP)
 - 405 subjects with schizophrenia (SZ) patients and normal controls (NC) with a single fMRI task
 - Additional subjects are relatives of SZ patients and intermediate disorder phenotypes
 - 32 demographic and clinical assessment measures



Demographics Overview

Table: Demographics for the PNC and BSNIP datasets.

PNC Dataset

	Males	Females	p-value
Age (years)	14.48 ± 3.32	14.69 ± 3.42	< 0.3042
European Anc. (EA)	316	303	< 0.5415
African Anc. (AA)	224	311	$< 10^{-5}$
WRAT Score	103.66 ± 16.56	101.45 ± 15.90	< 0.0212
	EA	AA	
Age (years)	14.59 ± 3.50	14.59 ± 3.23	< 0.9994
WRAT Score	108.68 ± 14.84	95.68 ± 14.80	$< 10^{-48}$

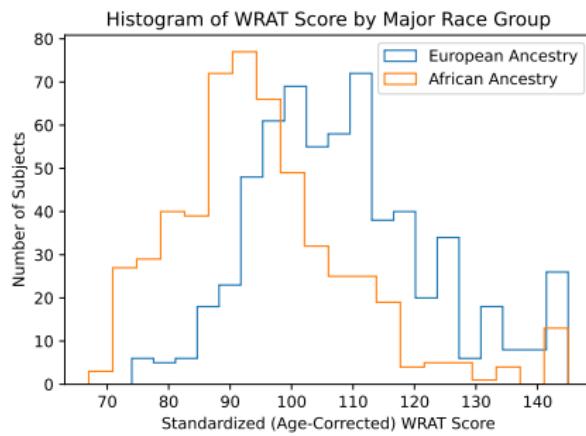
BSNIP Dataset

	Males	Females	p-value
Age (years)	35.24 ± 11.98	38.93 ± 12.45	< 0.0026
Caucasian Anc. (CA)	139	109	< 0.0221
African Anc. (AA)	82	75	< 0.5343
SZ Diagnosis	130	55	$< 10^{-10}$



Example of Confounding Effect on WRAT Score

- One of the 169 PNC fields is the wide range achievement test (WRAT) score
- Available raw and adjusted for age
- No adjustment for ethnicity, gives opportunity for confounding effects



Experiments

- We perform five experiments on the PNC and BSNIP datasets to investigate the ability of the DemoVAE model
 - To remove the confounding effects of demographics
 - To generate high quality synthetic data
 - To alter fMRI based on changing subject demographics or scanner task
- ① Prediction of WRAT Score Using DemoVAE Latents
 - ② Validation of fMRI Samples Generated by DemoVAE
 - ③ Phenotype Prediction Using DemoVAE Synthetic Data
 - ④ Correlation of Clinical Measures with DemoVAE Latents
 - ⑤ Imputation of fMRI Scanner Task



Prediction of WRAT Score Using DemoVAE Latents

Table: RMSEs (mean and standard deviation) of predicting standardized WRAT scores using fMRI FC input, SNP input, DemoVAE fMRI latents, DemoVAE SNP latents, and scalar race variable.

Input	WRAT Prediction RMSE
Null Model	15.18
Race Only	13.91 ± 4.65
Rest FC	14.73 ± 6.30
Nback FC	14.44 ± 6.77
Emoid FC	14.46 ± 7.10
SNPs	14.03 ± 7.35
Rest DemoVAE Latents	15.20 ± 0.27
Nback DemoVAE Latents	15.18 ± 0.23
Emoid DemoVAE Latents	15.18 ± 0.26
SNP DemoVAE Latents	15.14 ± 2.25



TULANE
UNIVERSITY

Validation of fMRI Samples Generated by DemoVAE

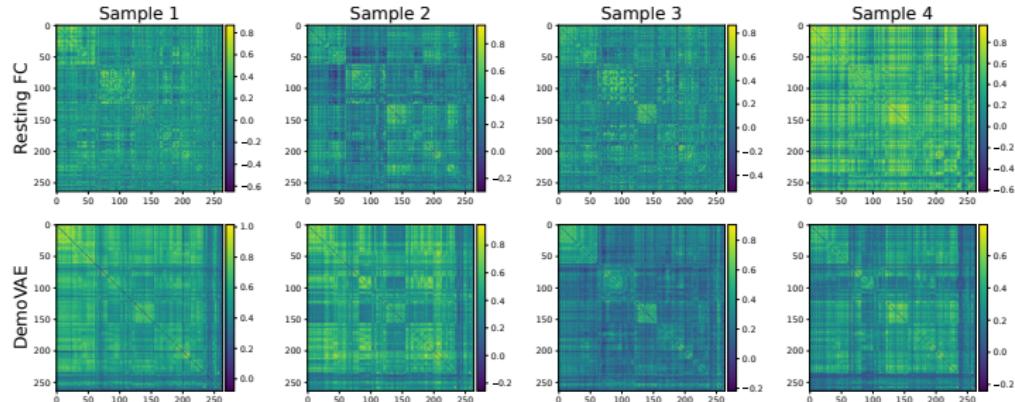


Figure: Sampled FC matrices for real PNC resting state scans (top) compared to synthetic DemoVAE (bottom). Qualitatively, DemoVAE generates convincing data.



t-SNE Embedding of DemoVAE Synthetic Data

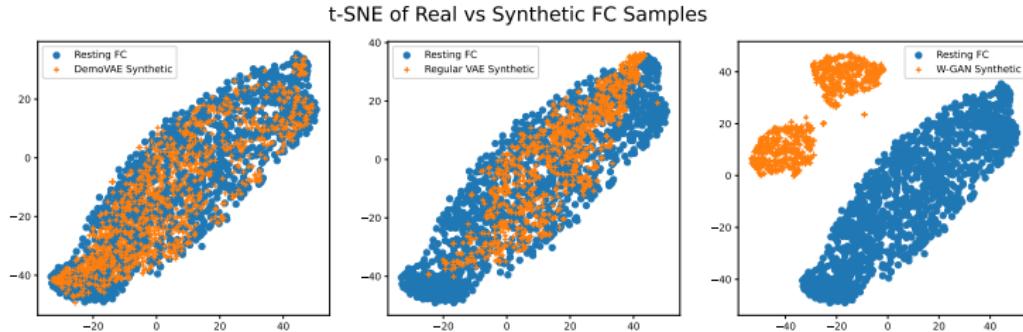


Figure: t-SNE embeddings of synthetic FC data from DemoVAE, traditional VAE, and W-GAN models overlayed on top of t-SNE embeddings of real resting state FC data from the PNC dataset. Blue circles represent embeddings of real subject FC data while orange crosses represent embeddings of synthetic data.



DemoVAE Synthetic Data Recapitulates Real Group Differences

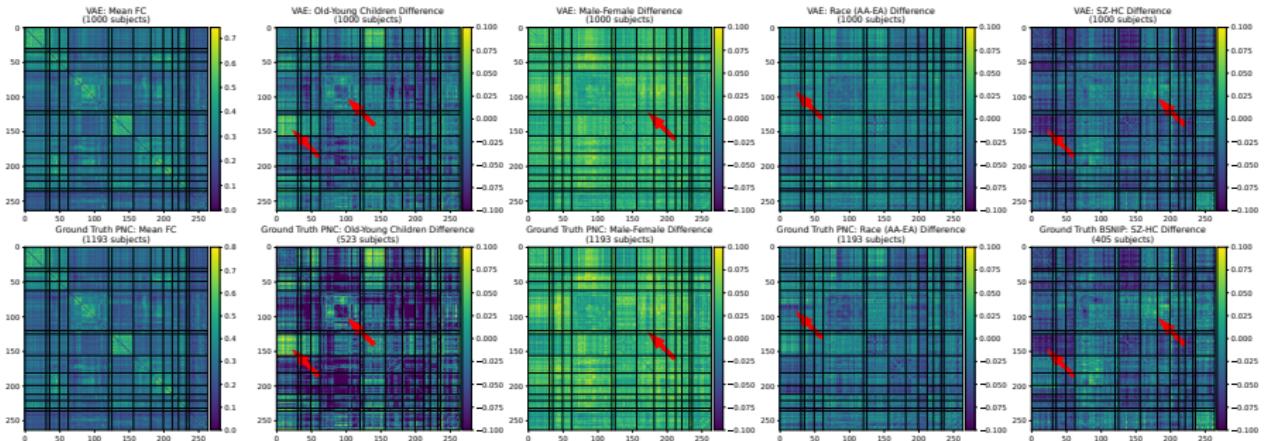


Figure: Group FC differences using real data and synthetic data generated by DemoVAE conditioned on appropriate demographic input. Top: synthetic DemoVAE data, bottom: real data. We see that DemoVAE qualitatively recapitulates group differences in the PNC and BSNIP datasets.

Phenotype Prediction Using DemoVAE Synthetic Data

Table: Transfer of models between fMRI and VAE. RMSE (age prediction) and mean accuracy (sex, race, and SZ prediction) for MLP models trained on ground truth fMRI data and tested on DemoVAE generated samples and vice versa.

Train on fMRI, Test on VAE
PNC Dataset

Predictive Task	Null Model	Rest FC	Rest PCFC	Nback FC	Nback PCFC	Emoid FC	Emoid PCFC
Age (years, RMSE)	3.30	0.570	2.181	0.468	1.97	0.495	1.91
Sex (ACC, %)	53.2	100	99.7	100	99.6	99.9	99.4
Race (ACC, %)	53.0	100	99.8	100	99.9	100	100

Train on VAE, Test on fMRI
PNC Dataset

Predictive Task	Null Model	Rest FC	Rest PCFC	Nback FC	Nback PCFC	Emoid FC	Emoid PCFC
Age (years, RMSE)	3.30	2.032	2.752	1.848	2.567	1.953	2.597
Sex (ACC, %)	53.2	88.9	90.4	90.9	91.2	91.1	91.7
Race (ACC, %)	53.0	93.2	96.3	93	96.1	93.1	97.1



Phenotype Prediction (BSNIP)

Table: Prediction cont'd (BSNIP dataset).

Train on fMRI, Test on VAE
BSNIP Dataset

Predictive Task	Null Model	FC	PCFC
Age (years, RMSE)	12.4	3.67	3.86
Sex (ACC, %)	54.5	100	100
Race (ACC, %)	61.2	100	100
Schizophrenia (ACC, %)	54.3	100	98.7

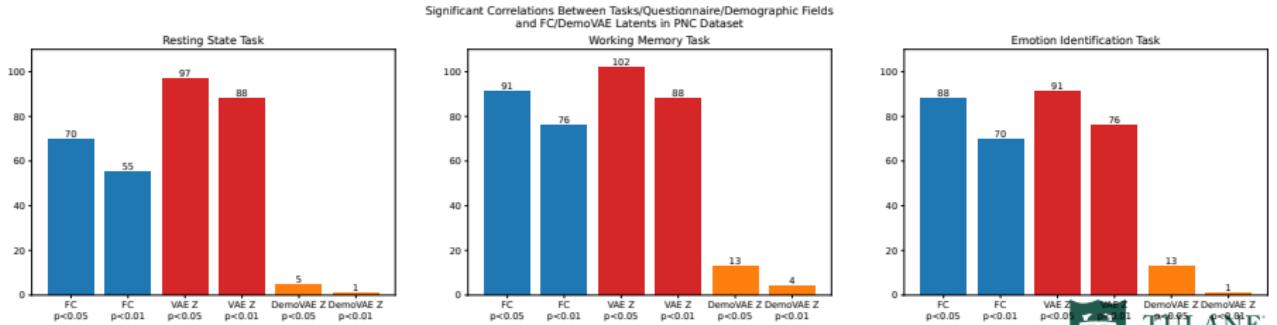
Train on VAE, Test on fMRI
BSNIP Dataset

Predictive Task	Null Model	FC	PCFC
Age (years, RMSE)	12.4	7.98	10.3
Sex (ACC, %)	54.5	97.5	94.5
Race (ACC, %)	61.2	96.0	93.3
Schizophrenia (ACC, %)	54.3	93.3	92.3



Correlation of Clinical Measures with DemoVAE Latents

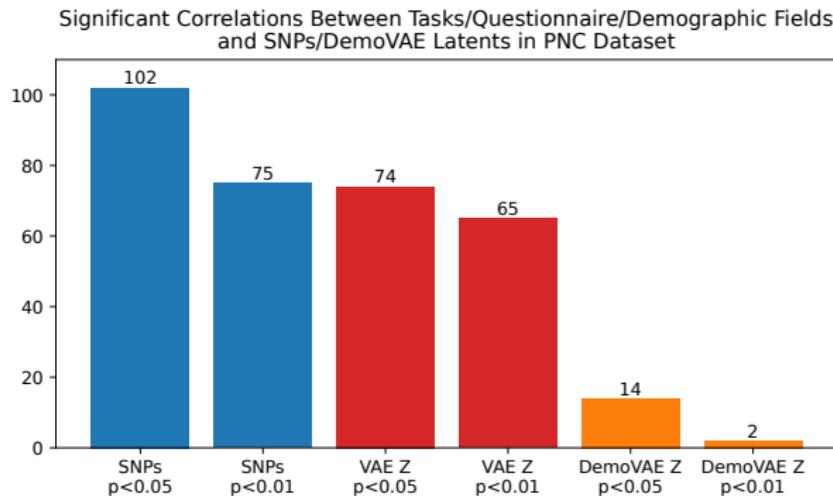
- 169 clinical questionnaire and computerized battery fields in PNC dataset
- 32 in BSNIP dataset
- More than half correlated with FC or regular VAE latents, almost no significant correlations with DemoVAE latents
- Blue = FC data, Red = VAE latents, Orange = DemoVAE latents



TULANE
UNIVERSITY

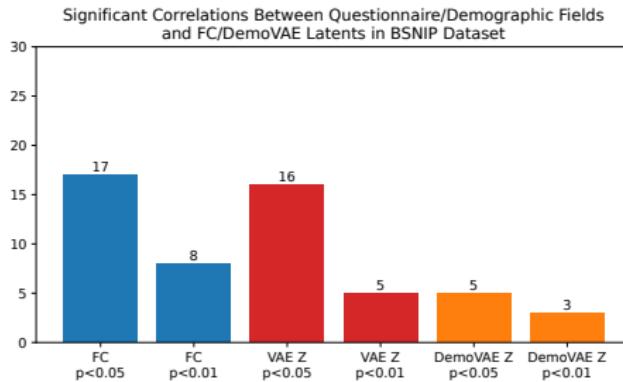
Correlations of Clinical Measures (PNC, SNPs)

- Similar trend as FC in genomic data (SNPs) for PNC dataset



Correlation of Clinical Measures (BSNIP)

- In BSNIP dataset, the fields having significant correlation with DemoVAE latents were:
 - total positive symptom score ($p < 0.0098$), total negative symptom score ($p < 0.0296$), total general symptom score ($p < 0.0011$), and total PANSS score ($p < 0.00033$)
 - taking or not taking anti-psychotics, $p < 0.0218$



Imputation of fMRI Scanner Task

- DemoVAE allows for superior imputation of scanner task by changing the scanner task field in decoder

Table: RMSEs (mean and standard deviation) for the reconstruction of one task FC from another scanner task in the test set, using MLP model, mean difference on training set, and DemoVAE. DemoVAE is used in deterministic mode and using best and average of 10 samples using 10% noise in the latent dimension.

	Rest → Nback	Rest → Emoid	Nback → Rest	Nback → Emoid	Emoid → Rest
Zero FC	0.333 ± 0.091	0.344 ± 0.100	0.368 ± 0.103	0.344 ± 0.100	0.368 ± 0.103
Reuse Input	0.218 ± 0.058	0.234 ± 0.064	0.232 ± 0.059	0.197 ± 0.054	0.234 ± 0.064
Average Training Set Diff.	0.210 ± 0.057	0.231 ± 0.064	0.230 ± 0.059	0.194 ± 0.055	0.231 ± 0.064
MLP	0.191 ± 0.047	0.200 ± 0.050	0.216 ± 0.047	0.190 ± 0.048	0.217 ± 0.050
DemoVAE (Deterministic)	0.196 ± 0.0526	0.202 ± 0.055	0.220 ± 0.056	0.188 ± 0.048	0.219 ± 0.059
DemoVAE (Best of 10)	0.176 ± 0.028	0.185 ± 0.034	0.203 ± 0.034	0.185 ± 0.032	0.203 ± 0.034
DemoVAE (Avg. of 10)	0.197 ± 0.052	0.202 ± 0.056	0.220 ± 0.056	0.189 ± 0.048	0.219 ± 0.059



Future Work

- We show that most FC-based predictive tasks are confounded by demographics (age, sex, race)
- Once demographics are controlled for in the latent dimension, most predictive ability disappears
- DemoVAE allows for imputation of scanner task as well as changing subject parameters (e.g., aging by 10 years, changing sex, etc.)
- Future work is to perform harmonization by using scanner site as a “demographic field,” similar to how we already use scanner task

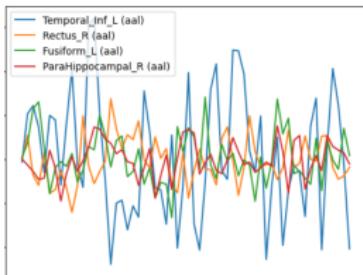


Questions

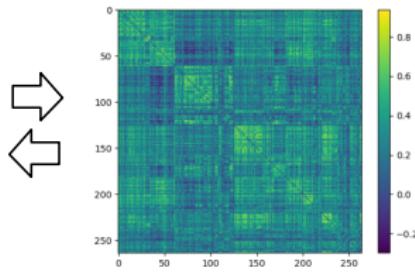
Thank you! Any questions?



Appendix: Generation of Timeseries



BOLD Timeseries



Functional Connectivity

- We train the VAE on FC features
- Is there a way to derive timeseries consistent with an FC matrix, so that we can derive alternate measures of connectivity?
 - E.g., such as partial correlation-based functional connectivity
- Yes, we can!

Output of DemoVAE Decoder

- FC is a positive semi-definite (PSD) matrix
- The output of the DemoVAE decoder $\hat{\mathbf{x}} = D_\theta(\mathbf{z}, \mathbf{y})$ can be either vectorized upper triangle of FC
 - In which case, it may have several negative eigenvalues
- Or a factor/factor transpose pair that reconstructs the FC matrix
 - Which may be rank-deficient, depending on the dimension of \mathbf{A}

$$\begin{aligned}\hat{\mathbf{X}}^{(1)} &= \hat{\mathbf{X}}_U + \hat{\mathbf{X}}_U^\top + \mathbf{I} \\ \hat{\mathbf{X}}^{(2)} &= \mathbf{A}\mathbf{A}^\top.\end{aligned}\tag{10}$$



Deriving the Covariance Matrix

- We can set any negative eigenvalues of the reconstructed FC to zero
- This doesn't seem to adversely effect predictive ability

$$\lambda_{\hat{\mathbf{X}},i} = \begin{cases} \lambda_{\hat{\mathbf{X}},i}, & \lambda_{\hat{\mathbf{X}},i} \geq 0 \\ \beta, & \text{otherwise} \end{cases} \quad (11)$$

- We then choose the standard deviation of timeseries $\sigma_i \in N(\mu_{\sigma_i}, \tau_{\sigma_i}^2)$
- And reconstruct the covariance matrix

$$\boldsymbol{\Sigma} = (\mathbf{1}\sigma^\top)\hat{\mathbf{X}}(\sigma\mathbf{1}^\top), \quad (12)$$



Cholesky Decomposition

- We then obtain a Cholesky-like decomposition, which is only available for positive definite (PD), via a QR decomposition and the eigenvalues/eigenvectors of Σ

$$\begin{aligned}\Sigma &= \mathbf{V} \Lambda \mathbf{V}^\top \\ &= (\mathbf{V} \sqrt{\Lambda})(\mathbf{V} \sqrt{\Lambda})^\top \\ &= (\mathbf{Q} \mathbf{R})^\top (\mathbf{Q} \mathbf{R}) \\ &= \mathbf{L} \mathbf{L}^\top.\end{aligned}\tag{13}$$

Timeseries via the Cholesky Factor

- Timeseries may then be constructed based on the property of the Cholesky decomposition that a standard normal random variable $X \in \mathcal{N}(0, 1)$ multiplied by the Cholesky factor \mathbf{L} creates a multivariate normal variable vector with zero mean and covariance matrix $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^\top$.
- This is compatible with fMRI data, which are usually bandpass filtered prior to analysis.
- It assumes, however, that the timeseries BOLD signal is stationary, which is sufficient when producing correlation-based metrics.

