

Prospectus

Analysis of Functional Connectivity and its Fusion with Multi-Omics
Using Contrastive Learning

Anton Orlichenko
Advisor: Yu-Ping Wang

June 2023



Outline

- 1 Background on fMRI and Cognitive Science
- 2 Problem Statement and Goals
- 3 Specific Aims
 - Aim 1: Latent Similarity for Small Sample Size, High Dimensionality Datasets
 - Aim 2: ImageNomer, A Data Exploration Tool Reveals Racial Confound in fMRI Data
 - Aim 3: Angle Basis: A Generative Model and Decomposition of Functional Connectivity
 - Aim 4: Contrastive Learning for Fusing Omics with Brain Imaging
- 4 Summary and Acknowledgements

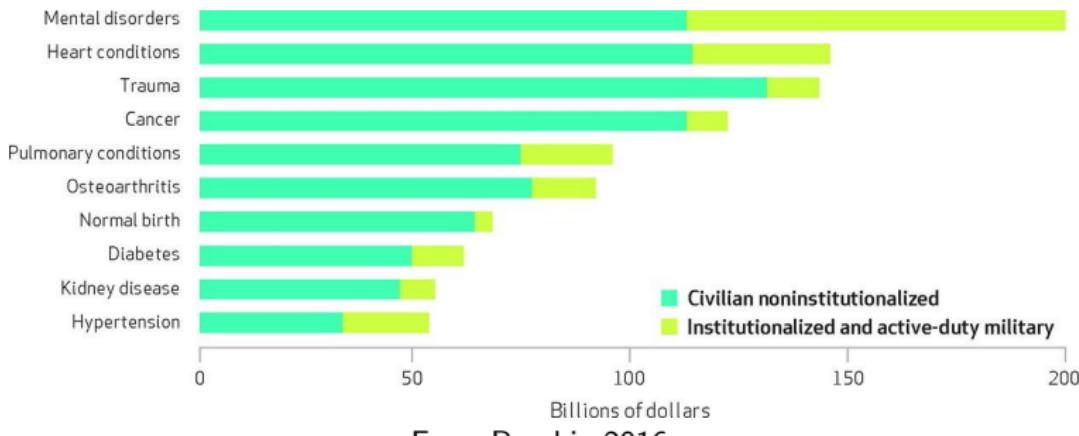


Background on fMRI and Cognitive Science



Clinical Problem

- Schizophrenia, ADHD, depression, and other mental illnesses cost the U.S. \$201+ billion annually¹
- Dementia and Alzheimer's cost the U.S. \$157+ billion annually²
- Diagnosis of these diseases may be unreliable until symptoms become severe, when treatment options are more limited



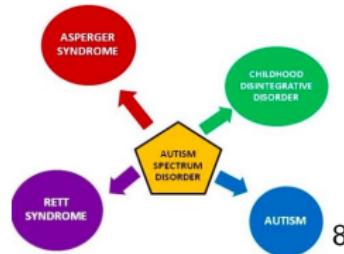
From Roerhig 2016.

¹ Roerhig 2016 <https://doi.org/10.1377/hlthaff.2015.1659>

² Hurd et al. 2013 doi:10.1056/NEJMsa1204629

Statistics

- 1 in 300 people are living with schizophrenia³
- 1 in 10 children may be affected by ADHD⁴
- 1 in 10 people may have had a major depressive episode in the past several years⁵
- 1 in 36 children may be diagnosed with autism spectrum disorder⁶
- 1 in 9 Americans over 65 years old are living with Alzheimer's⁷



³ Desai et al. 2013 10.1111/jphs.12027

⁴ Dincer et al. 2022 10.1177/10870547221099963

⁵ Brody et al. 2018 Prevalence of depression among adults aged 20 and over

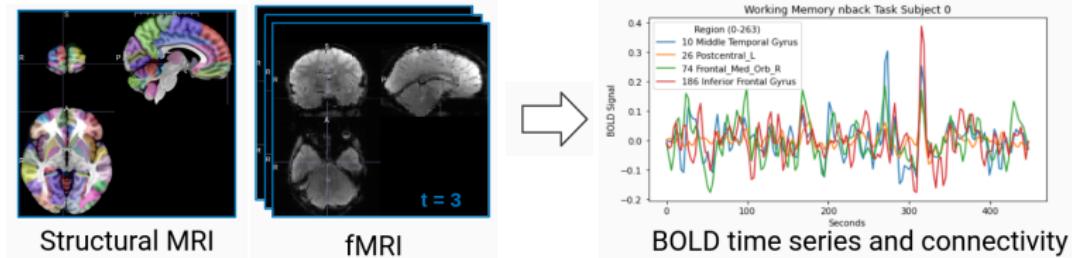
⁶ Maenner et al. 2020 10.15585/mmwr.ss7011a1

⁷ 2023 Alzheimer's disease facts and figures

⁸ <https://speechandot.com/what-are-the-types-of-autism-spectrum-disorder/>

fMRI

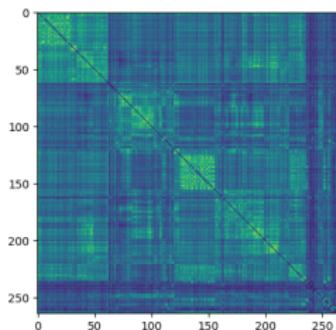
- Functional magnetic resonance imaging (fMRI) provides a non-invasive estimate of brain activity by exploiting the blood oxygen level-dependent (BOLD) signal



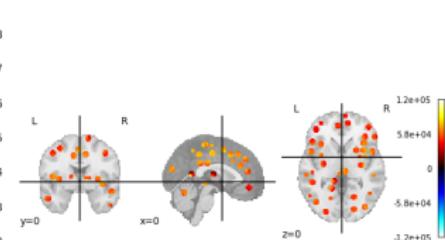
Example extraction of BOLD timeseries from fMRI volumes.

Functional Connectivity (FC)

- Functional connectivity (FC) is the temporal Pearson correlation of BOLD signal between brain regions
- These brain regions are typically found using either ICA⁹ or an atlas¹⁰
- FC or a variant is the most common starting point for predictive fMRI studies



mean FC of subjects in dataset



using 264 region Power template

⁹Calhoun et al. 10.1016/j.neuroimage.2008.10.057

¹⁰Power et al. 10.1016/j.neuron.2011.09.006

Power Atlas Functional Networks

- Example of Power atlas partition of brain into functional networks

Functional Networks

ROIs		ROIs	
0-29	Somatomotor Hand	156-180	Frontoparietal
30-34	Somatomotor Mouth	181-198	Salience
35-48	Cinguloopercular	199-211	Subcortical
49-61	Auditory	212-220	Ventral Attention
62-119	Default Mode	221-231	Dorsal Attention
120-124	Memory	232-235	Cerebellar
125-155	Visual	236-263	Uncertain



fMRI and Mental Health

- fMRI has been used to predict disease status and (endo)phenotypes such as age, sex, and general fluid intelligence¹¹
- Machine learning predictions of brain age have been correlated with future Alzheimer's diagnosis¹²
- Most diagnoses of mental disorders are still made by psychiatrists or physicians based on cognitive tests¹³

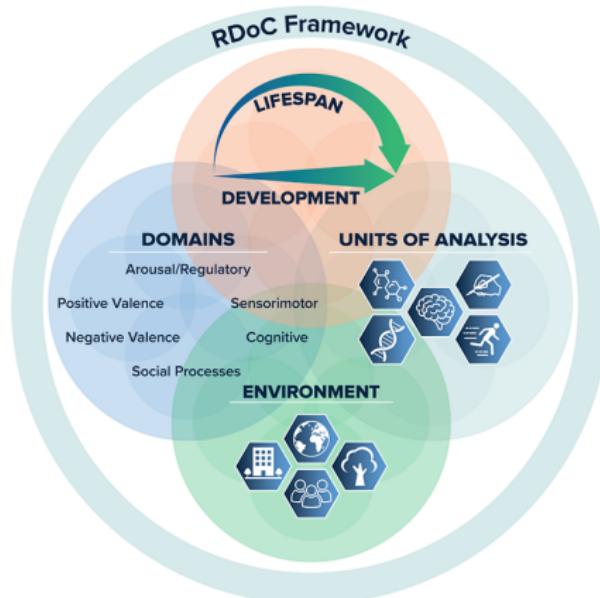
¹¹ Qu et al. 2021 10.1109/TBME.2021.3077875

¹² Millar et al. 2022 10.1016/j.neuroimage.2022.119228

¹³ <https://www.ndcn.ox.ac.uk/divisions/fmrib/what-is-fmri/how-is-fmri-used>

Research Domain Criteria (RDoc)

- The National Institutes of Mental Health (NIMH) have created the Research Domain Criteria (RDoC) to put mental health diagnoses on a more rigorous scientific basis¹⁴



15

¹⁴ Insel et al. 10.1176/appi.ajp.2010.09091379

¹⁵ <https://www.nimh.nih.gov/research/research-funded-by-nimh/rdoc/about-rdoc>

fMRI in the Clinic

- Clinically, fMRI is used for pre-surgical planning and to follow the progression of concussion¹⁶



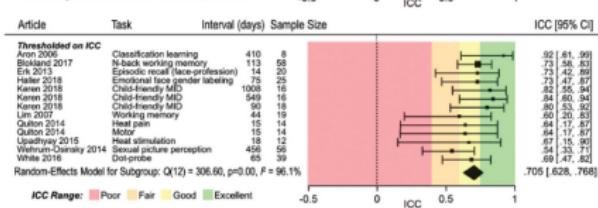
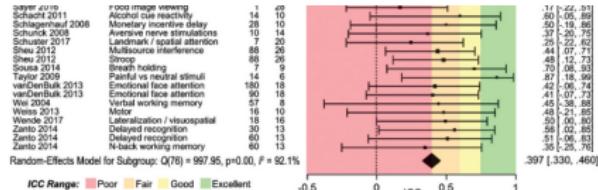
Example of a clinical lab setup for monitoring concussion¹⁷

¹⁶ Kaushal et al. 2019 10.1002/hbm.24440

¹⁷ <https://www.cognitivefxusa.com/blog/fmri-brain-scans-duke-study-implications>

Potential Pitfalls

- Voxel-based fMRI analysis has been criticized for having poor test-retest replicability¹⁸
- This may explain why FC has become popular as a predictive tool



¹⁸ Elliott et al. 2020 10.1177/0956797620916786

Problem Statement and Goals



Problem Statement

We seek to better understand the structure and function of the human brain, in order to

- ① characterize neurological disease for diagnosis and treatment
- ② understand normal development

but come up against **two challenges and one opportunity.**



Challenges

High Cost of Data Acquisition

The cost of acquiring an fMRI scan for a single subject is \$500-\$1000.^a

^aSzucs and Ioannidis 2020 10.1016/j.neuroimage.2020.117164

- Same magnitude of cost for acquiring genomic data.

Small Sample Size, High Dimensionality

The median cohort size for fMRI studies in 2017 and 2018 was 23 subjects,^a while

- number of voxel-level features can be $> 10^6$
- number of connectivity-level features can be $> 10^4$

^aSzucs and Ioannidis 2020 10.1016/j.neuroimage.2020.117164

Opportunity

Large Amount of Unlabeled Data

Repositories exist for large amounts of unlabeled or out-of-distribution data

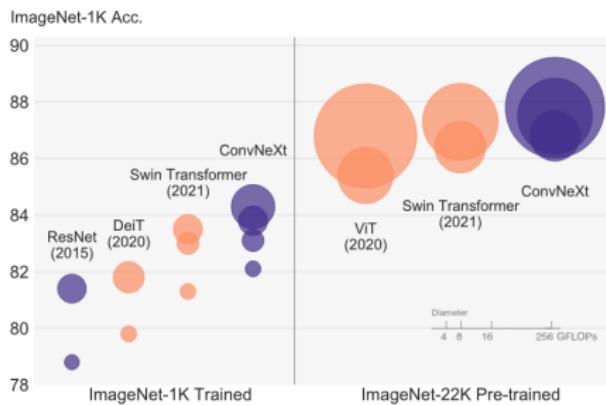
- UKB has 50,000 subjects with fMRI scans from a largely homogenous population^a
- OpenNeuro has 800+ small open access datasets with varied demographics^b

^aSudlow et al. 2015 10.1371/journal.pmed.1001779

^bMarkiewicz et al. 2021 10.7554/eLife.71774

Idea: Pre-training Improves Performance

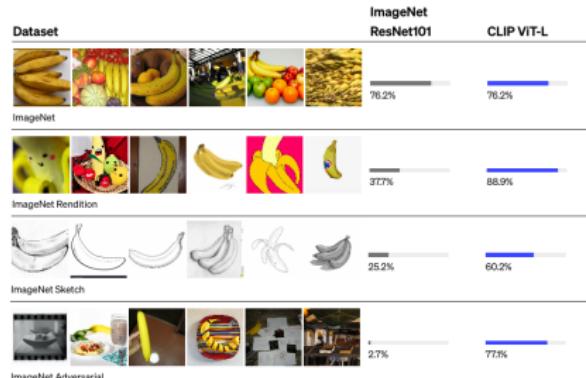
- It has been shown in the image classification and natural language processing AI literature that pre-training models on large unlabeled datasets improves performance



Pre-training yields a 4-6% improvement in image classification.¹⁹

Contrastive Learning Fuses Modalities

- Pre-training learns important features and often falls under the umbrella of **contrastive learning**
- CLIP fuses features of natural language and images, allowing, e.g., AI-generate art



Contrastive Language Image Pretraining.²⁰



²⁰ <https://openai.com/research/clip>

Goals

The goals of this dissertation research are the following:

- ① **Aim 1: Improve model performance** in the
small sample size, high dimensionality regime
- ② **Aim 2: Create tools** for data exploration, quality control, and
detection of unreported covariates/confounders
- ③ **Aim 3: Construct a generative model** for the decomposition of FC
to create augmentations for use in pre-training
- ④ **Aim 4: Find associations** between brain networks and genotype or
gene expression using contrastive learning



Specific Aims



Aim 1: Latent Similarity for Small Sample Size, High Dimensionality Datasets



Problem of High Dimensionality

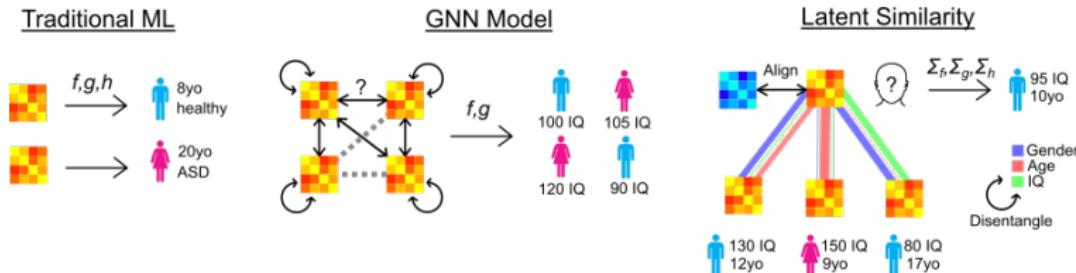
Question

Can we create a model that is robust with very low sample sizes and high dimensional features?

- Makes very conservative predictions with a small training set
- Increases in predictive power with a large training set



Solution: Latent Similarity



- As in contrastive learning, we explicitly model subject-subject latent similarities²¹²²
- Traditional ML models (Ridge, MLP, CNN) do not explicitly consider inter-subject relationships
- GNN models require many degrees of freedom to estimate edges

²¹ Orlichenko et al. 2022 10.1109/TBME.2022.3232964

²² Tool available at <https://github.com/aorliche/LatentSimilarity>

Similarity Kernel

- We model a linear similarity kernel to project data to a very low dimensional subspace
- Followed by softmax aggregation
- Can perform cross-modality alignment and disentanglement
- A similar approach was recently used in the image domain²³

$$\text{sim}(a, b) = \langle \phi(\mathbf{x}_a), \phi(\mathbf{x}_b) \rangle$$

$$= \mathbf{x}_a \mathbf{A} \mathbf{A}^T \mathbf{x}_b^T,$$

$$\mathbf{E} = S_{Row}((\mathbf{1} - \text{diag}(\infty)) \odot \mathbf{X} \mathbf{A} \mathbf{A}^T \mathbf{X}^T), \quad (1)$$

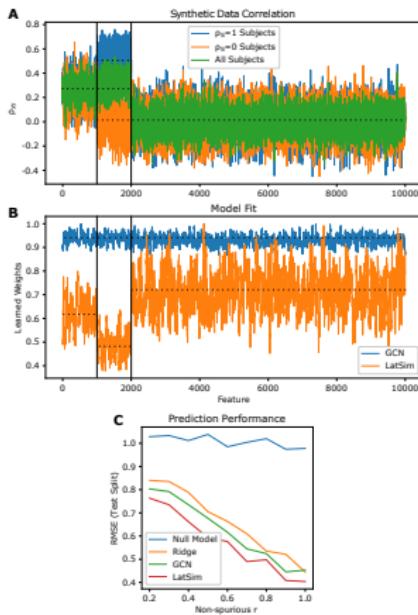
$$S(\mathbf{z})_i = \frac{e^{z_i/\tau}}{\sum_{j=0}^N e^{z_j/\tau}},$$



²³ Zheng et al. CVPR 2022 10.48550/arXiv.2203.06915

Simulation Results

- We find that LatSim performs better in simulation studies on synthetic data

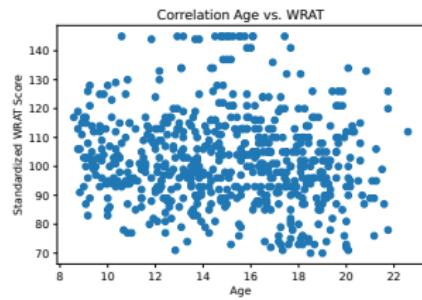
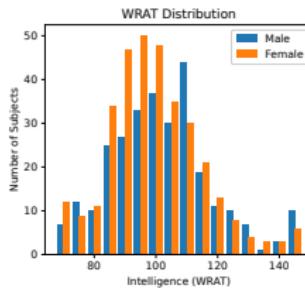
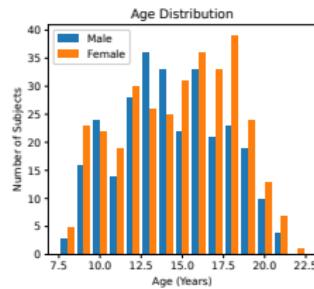


PNC Dataset

- We test on Philadelphia Neurodevelopmental Cohort (PNC), using reduced number of subjects

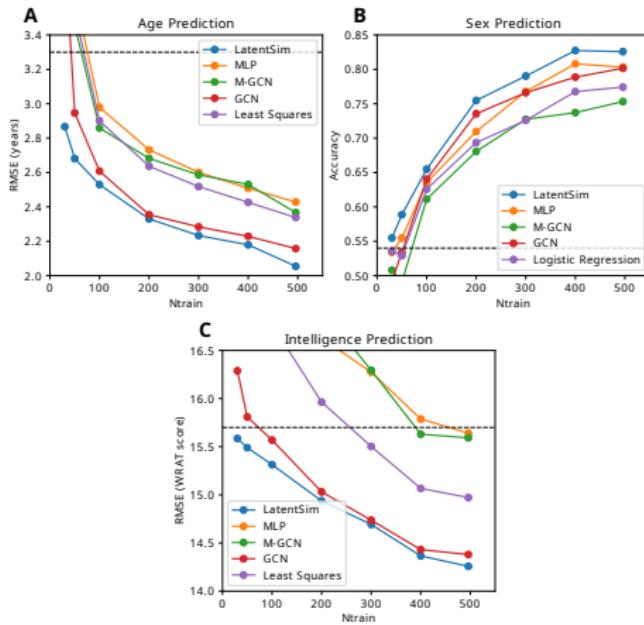
	Number of Subjects
Males	286
Females	334
Total	620

	Min	Mean	Max
Age (months)	103	180±39	271
Age (years)	8.6	15±3.3	22.6
WRAT score	70	102±15.7	145



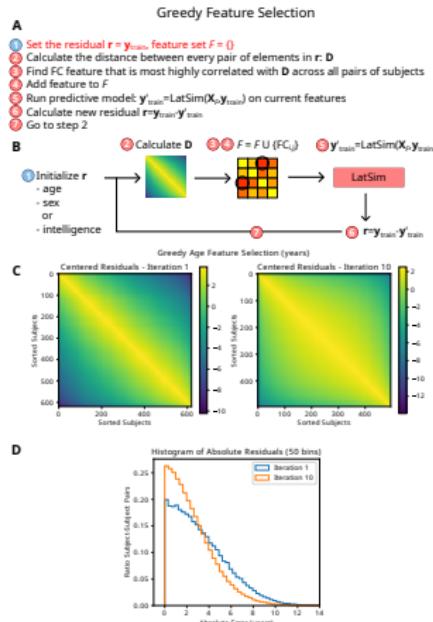
Latent Similarity FC Prediction

- We achieve superior predictive accuracy, especially at small sample sizes (predict age, sex, and WRAT score)



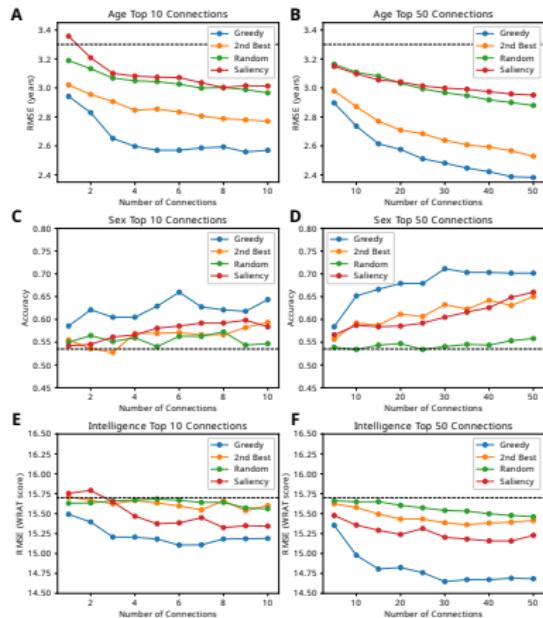
Latent Similarity Greedy Selection

- We include a greedy feature selection algorithm which uses LatSim in its inner loop



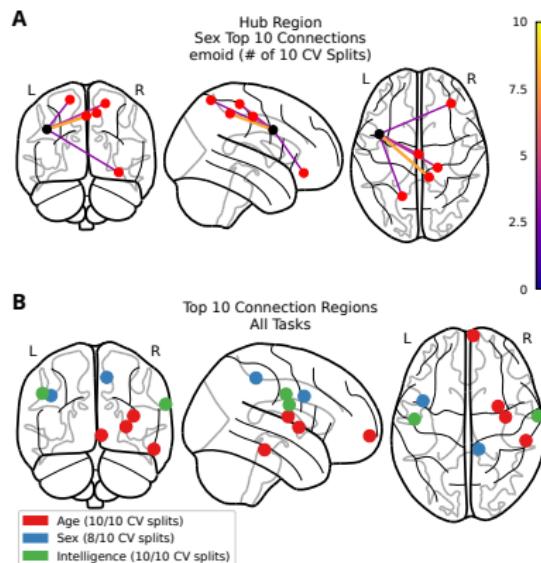
Superior Performance of Greedy Selection

- We found bootstrap greedy selection features to be superior to 2nd choice, random features, and saliency (gradient-based method)



Key Connections

- We identified several key connections useful for each predictive task
- However, these connections may be replaced by others with small to medium predictive penalty



Computational Benefits

- LatSim is much faster than all but Ridge or Least Squares models
- Why choose a slow model to get similar or worse results?

Table: Training time for all 10 folds of 10-fold cross validation.

Model	LatSim	Lstsq	Logistic	GCN	MLP	M-GCN
Epochs	200	-	100	1e4	1e4	5e3
Training Time	4.3s	<1s	63.4s	406s	364s	5912s



Latent Similarity Takeaways

Proof of Concept for Similarity Method

- We show that similarity-based methods, at the core of contrastive learning and data fusion, can achieve superior results at small sample sizes
- LatSim also competitive with other models at larger sample sizes
- Much faster than, e.g., scikit-learn implementation of logistic regression



Aim 2: ImageNomer, A Data Exploration Tool Reveals Racial Confound in fMRI Data



Problem of Dataset Exploration and Quality Control

Question

- How to quickly browse through and validate FC and genomic data?
- How to find correlations in phenotypes and visualize importance of features for prediction?

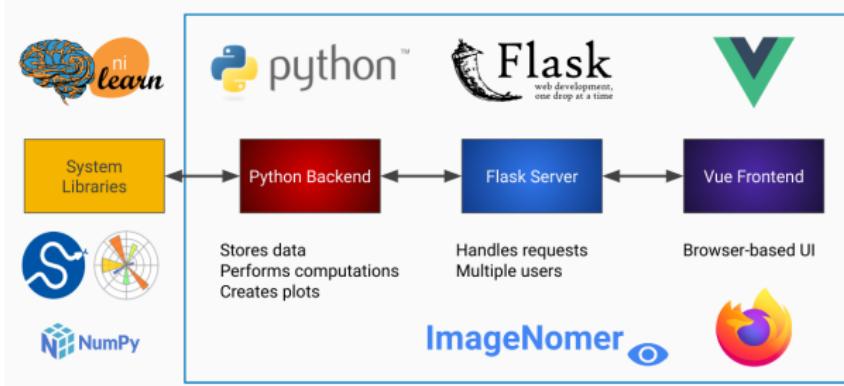
Preview of Findings

- We create an easy-to-use tool for exploring and visualizing correlations in an fMRI dataset
- We find a previously ignored confound in fMRI data affecting achievement score prediction



Our Solution

- We create a browser-based tool called ImageNomer to visualize subjects and get statistics about a dataset

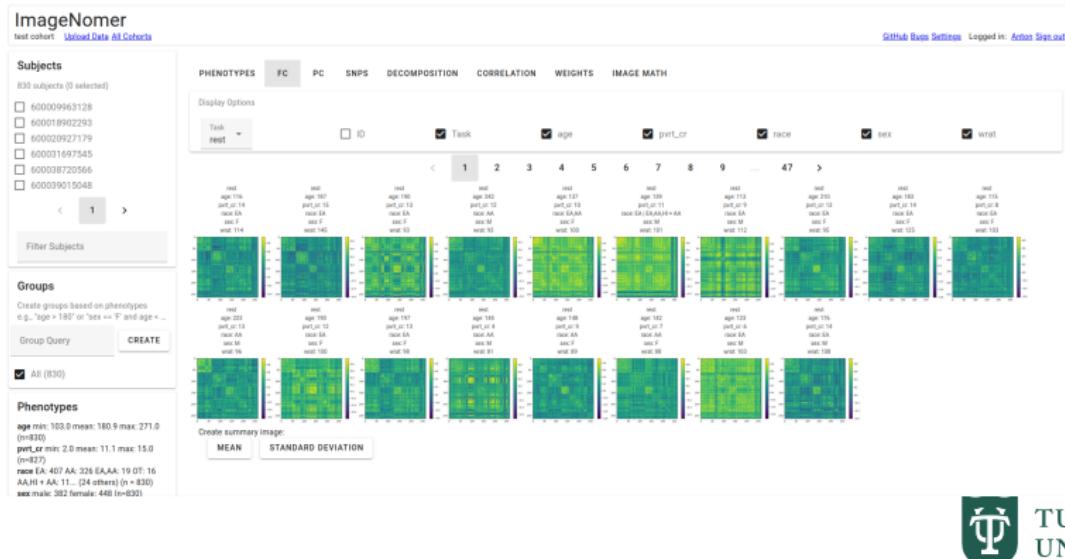


Architecture of the tool.²⁴

²⁴ Orlichenko et al. 2023 10.11117/12.2654311

ImageNomer Main View

- A screenshot of the main FC view in ImageNomer.
- Tool is available at ²⁵. See links for tutorial.
- <https://aorliche.github.io/ImageNomer/live/>



²⁵ <https://github.com/TulaneMBB/ImageNomer>

ImageNomer Documentation

- On-line documentation available at ReadTheDocs²⁶
- Tutorial walks through how to use features of ImageNomer
 - ▶ Phenotype distribution
 - ▶ Phenotype-phenotype or phenotype-FC correlation
 - ▶ Visualization of model weights
 - ▶ Image math

The screenshot shows two main sections of the ImageNomer documentation:

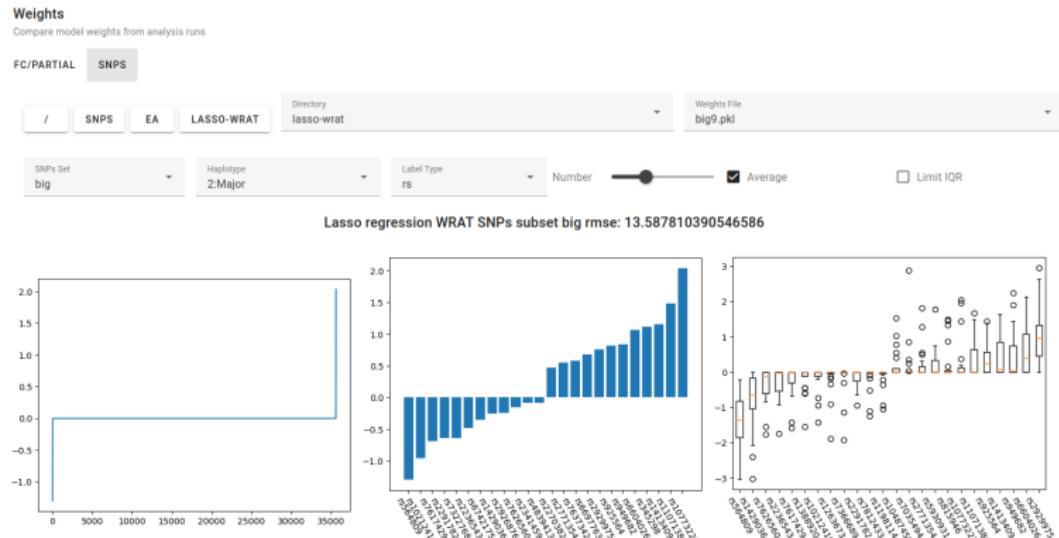
- Main Page:** Displays the "ImageNomer" logo, a brief description ("MRI and cancer viewer for machine learning"), and a "Goals" section listing four items: 1. Exploration of MRI data, 2. Quality control, 3. Correlation analysis, and 4. Visualization of MRI, model weights and distribution. It also includes a "Quick Start with Docker" section and a "Importing Data" section.
- Tutorial Section:** Titled "Fibromyalgia Dataset Tutorial", it provides instructions for opening the dataset in ImageNomer and navigating through various features like "Phenotype Correlation", "FC View", and "Suggested Workflows". A "Python application" section is also present.

26

<https://imagenomer.readthedocs.io/en/latest/index.html>

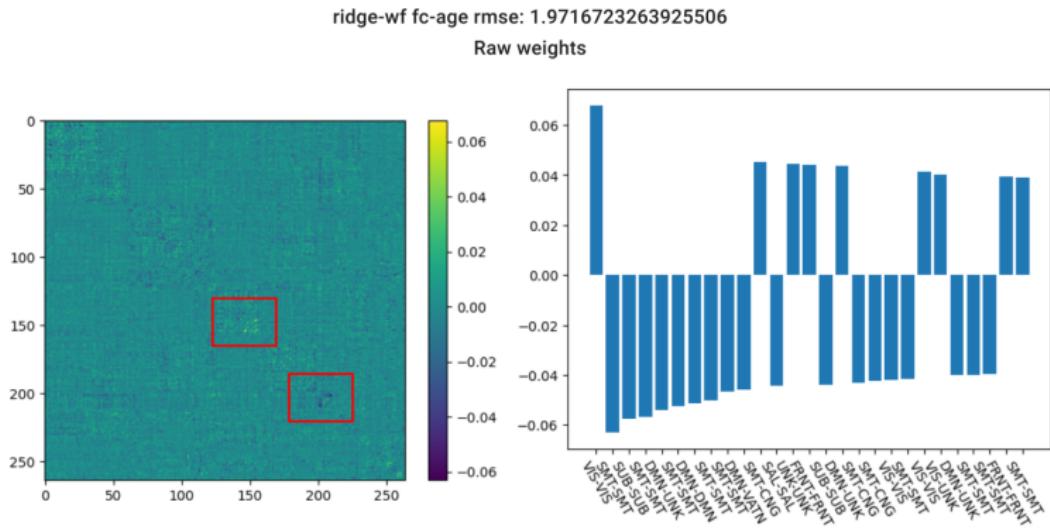
ImageNomer SNP Visualization

- Besides fMRI data, we can process SNP information



ImageNomer FC Model Weights Visualization

- Using the ImageNomer GUI, one can identify which regions machine learning models flag as important for prediction



Using Phenotype Explorer, Potential Race Confound is Apparent

- By creating plots of phenotype correlations, we see race or SES is a potential confound in WRAT score²⁷ prediction
- ...as long as FC data can be used to predict race

Group Selection Panel

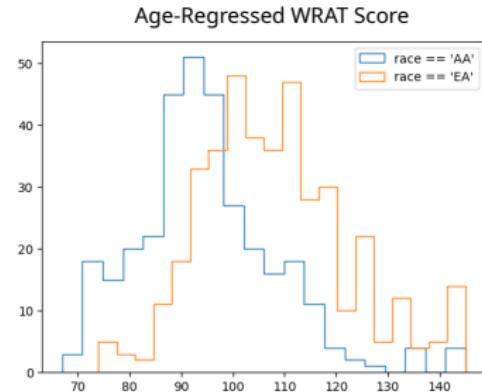
Groups

Create groups based on phenotypes
e.g., "age > 180" or "sex == 'F'" and age < ...

Group Query
`sex == 'F'`

CREATE

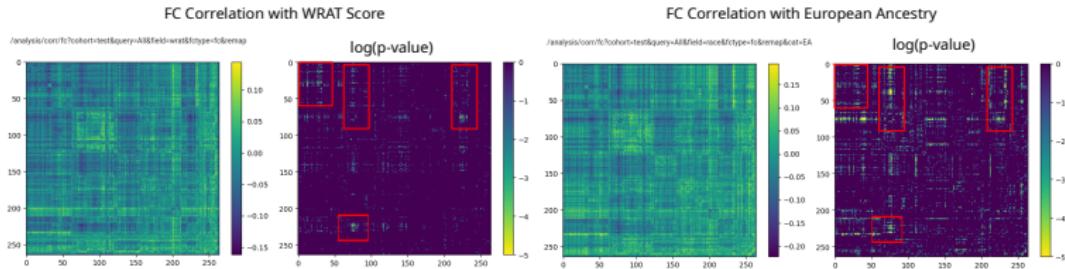
All (830)
 age < 160 (265)
 age > 160 (559)
 race == 'AA' (326)
 race == 'EA' (407)
 sex == 'M' (382)
 sex == 'F' (448)



²⁷ Sayegh et al. 2014 10.1093/arclin/acu059

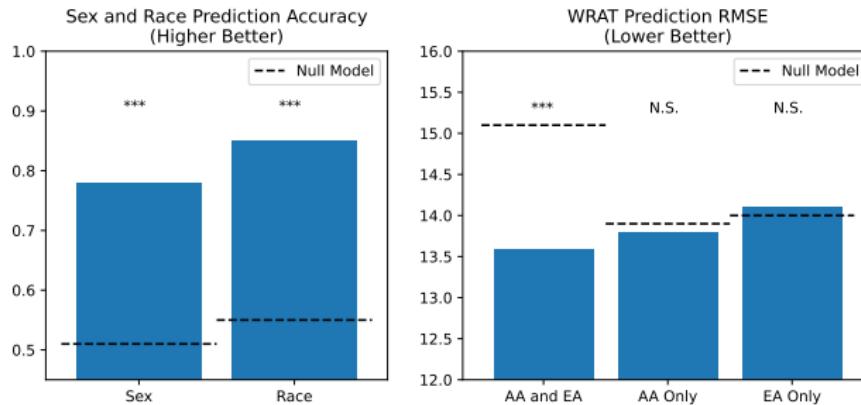
Achievement-FC Correlation is Subset of Race-FC Correlation

- We perform correlation analysis between FC and WRAT score and apply Bonferroni-corrected t-test
- We find achievement-FC correlation is a subset of race-FC correlation



Race Can Be Predicted, Unbiased Achievement is Harder

- We test the effect of the confound by predicting WRAT score from FC on whole cohort and within race groups
- When controlling for race, WRAT score prediction not significantly better than guess



ImageNomer Takeaways

Proof of Concept for Large Dataset and Multi-Omics Work

- We create an omics and fMRI visualization tool that requires minimal programming experience to use
- We show that we are able to analyze datasets with
 - ▶ thousands of subjects
 - ▶ hundreds of demographic features
 - ▶ very high FC and SNP dimensionality
- In the process, we find a confound that was unaddressed in and potentially contaminated previous work



Aim 3: Angle Basis: A Generative Model and Decomposition of Functional Connectivity



Problem of fMRI Data for Contrastive Learning

Question

Contrastive learning typically uses augmentations or different views of images/text to create positive/negative pairs

- How to create data-driven augmentations for fMRI data?

Preview of Findings

- We create a theoretically-motivated generative model for FC
- Model has several uses besides augmentations



Angle Basis Overview

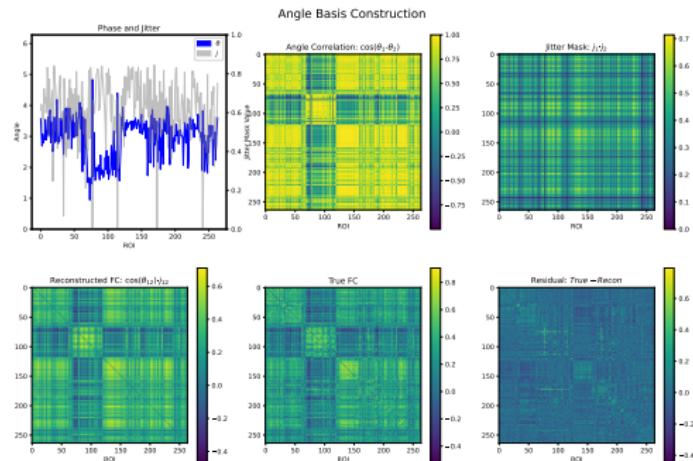
Key Features

- ① 10x data compression
- ② 97.3% fingerprinting identification (residual) vs 62.5% for FC
- ③ Modest 5% AUC prediction improvement over FC
- ④ Data augmentation for contrastive learning
- ⑤ Generation of synthetic FC from user-input patient characteristics
(network-level explainability/interpretability)
- ⑥ Does not rely on knowledge of a population; a single subject is enough



Example of Angle Basis Decomposition

- In searching for network level explanations for FC-based prediction, we created a simple yet powerful decomposition and compression of FC²⁸



Tool available at ²⁹

²⁸ Orlichenko et al. 2023 10.48550/arXiv.2305.10541

²⁹ <https://github.com/aorliche/AngleBasis>

Angle Basis Decomposition

- Based on the well-known phase lock value (PLV)

$$\begin{aligned}\mathbb{H}[x](t) &= \frac{1}{\pi} \text{p.v.} \int_{-\infty}^{\infty} \frac{x(\tau)}{t - \tau} d\tau \\ x_a(t) &= x(t) + i\mathbb{H}[x(t)] \\ x_a(t) &= a(t)e^{i\theta(t)} \\ \theta_{cd}(t) &= \theta_c(t) - \theta_d(t)\end{aligned}\tag{2}$$

$$\text{PLV}_{cd} = \frac{1}{T} \left| \sum_{t=1}^T e^{i\theta_{cd}(t)} \right|$$

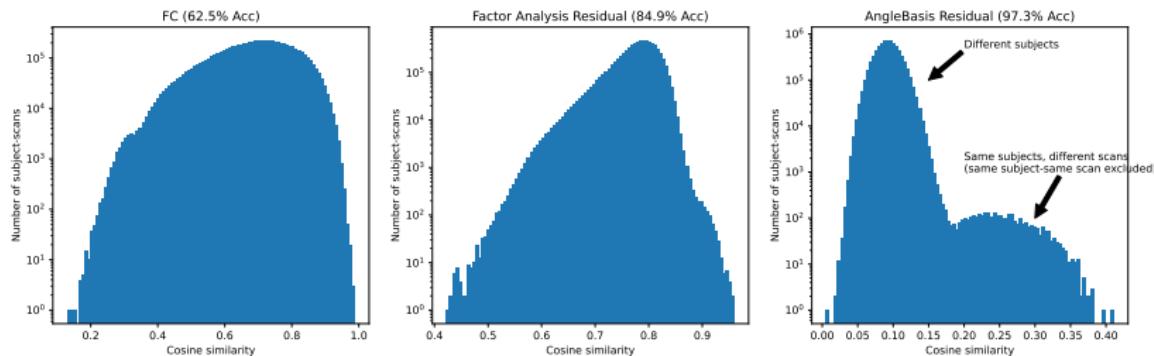
- We use a basis of orthogonal sine waves and jitter mask to compress and reconstruct FC

$$\begin{aligned}\hat{\rho}_{cd}^{(n)} &= j_c^{(n)} \cdot j_d^{(n)} \cdot \cos(\theta_c^{(n)} - \theta_d^{(n)}) \\ \tilde{\rho}_{cd} &= \frac{1}{N} \sum_{n=1}^N \hat{\rho}_{cd}^{(n)}\end{aligned}\tag{3}$$



Fingerprinting Accuracy

- We find that the residual of angle basis has 97.3% accuracy in identifying the same subject from different scans
- Compared to 62.5% for FC



Classification Improvement

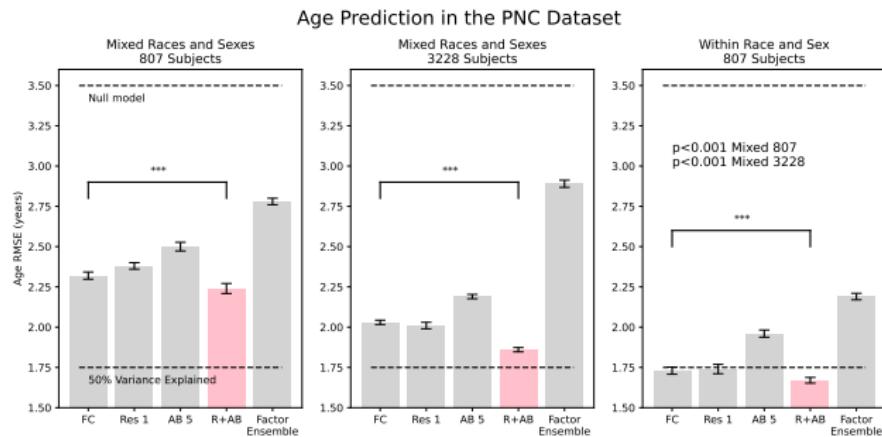
- On their own, angle basis and residual have equal predictive ability compared to FC
- Ensemble of the two is significantly better than FC

Classification Accuracy (AUC)					
Dataset	Predictive Task	FC	Deep AE Ens	AB+Res Ens	p-value
BSNIP	SZ/NC	0.785 ± 0.048	0.701 ± 0.062	0.804 ± 0.040	0.240
BSNIP	Sex	0.755 ± 0.023	0.645 ± 0.085	0.791 ± 0.022	0.002
BSNIP	Race	0.845 ± 0.022	0.739 ± 0.051	0.866 ± 0.023	0.108
PNC	Sex	0.886 ± 0.006	0.744 ± 0.065	0.923 ± 0.010	< 0.001
PNC	Race	0.946 ± 0.007	0.812 ± 0.040	0.973 ± 0.003	< 0.001
BSNIP→PNC	Sex	0.667 ± 0.017	0.601 ± 0.032	0.700 ± 0.013	< 0.001
BSNIP→PNC	Race	0.807 ± 0.018	0.710 ± 0.010	0.847 ± 0.012	< 0.001
PNC→BSNIP	Sex	0.629 ± 0.019	0.572 ± 0.019	0.667 ± 0.013	< 0.001
PNC→BSNIP	Race	0.800 ± 0.010	0.702 ± 0.022	0.832 ± 0.009	< 0.001

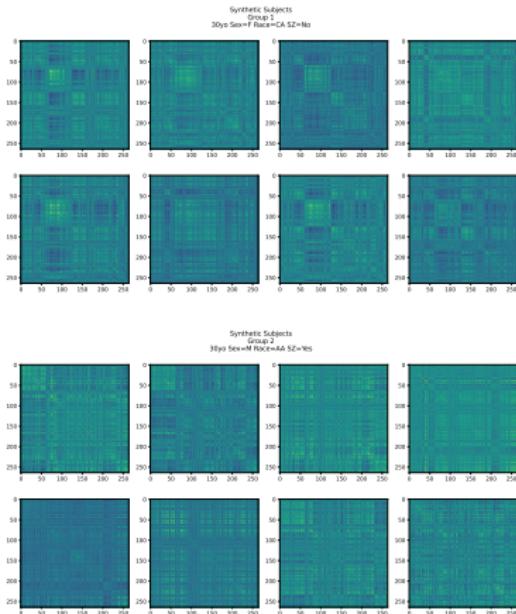


Importance of Within-Group Predictions

- Age regression results also show better prediction error with ensemble of angle basis and residual compared to just FC
- Moreover, we see the importance of making prediction within sex and race-matched groups



FC Generation for Explainability

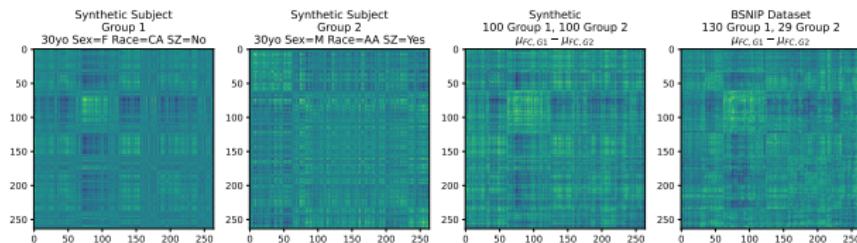


First 8 synthetic subjects from 30yo normal CA female group and 30yo schizophrenic AA male group.



Validation on Group Differences

- Synthetic subjects have similar group differences compared to demographic-matched real subject FCs
- Synthetic generation allows for extrapolation as well as visualization



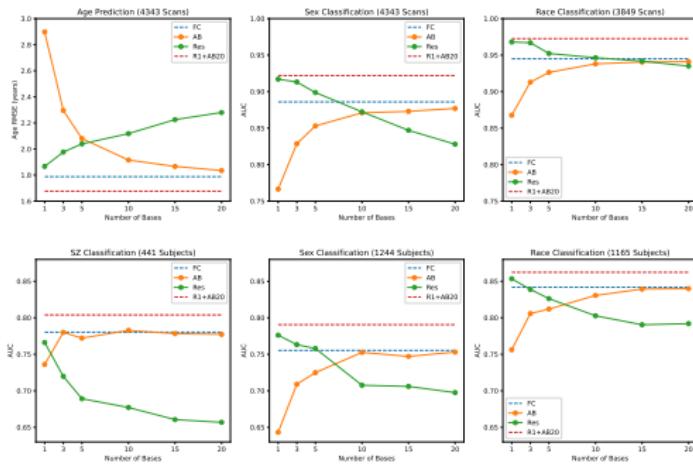
- Females are known to have increased intra-DMN connectivity³⁰



³⁰ Ficek-Tani et al. 2022 10.1093/cercor/bhac491

Prediction Accuracy as a Function of Number of Bases

- Both angle basis and residual have predictive value
- Varies inversely as the number of sinusoidal bases and jitter masks is increased



AngleBasis Takeaways

Proof of Concept for Augmentation and Manipulation of FC

- We are able to decompose, compress, and generate FC
 - ▶ to identify subjects with 97.3% accuracy
 - ▶ to achieve a modest 5% classification AUC gain
 - ▶ to generate synthetic FC based on patient demographics, increasing intuition and explainability
- In the process, we are now able to perform theoretically-motivated augmentations, as in the mixup model^a

^aZhang et al. 2017 10.48550/arXiv.1710.09412

Aim 4: Contrastive Learning for Fusing Omics with Brain Imaging



Problem of Leveraging Unlabeled Datasets for Discovering Causal Associations with Omics

Question

Can we use contrastive learning

- to leverage the large amounts of out-of-distribution data (UKB, OpenNeuro) available to create better models?
- to fuse brain imaging with multi-omics to discover links between omics and fMRI endophenotype?

Important

Work is only beginning on this Aim, thus many of the ideas in this section are speculative.



Foundations of Contrastive Learning

- The basic foundation of contrastive learning is an objective function based on positive and negative pairs
- The InfoNCE (noise contrastive estimation) loss is shown below:

$$\mathcal{L}_{CL} = -\frac{1}{N} \sum_i^N \log \frac{e^{\mathbf{q}_i^T \mathbf{k}_i^+ / \tau}}{e^{\mathbf{q}_i^T \mathbf{k}_i^+ / \tau} + \sum_{j=1}^K e^{\mathbf{q}_i^T \mathbf{k}_{i,j}^- / \tau}} \quad (4)$$

- \mathbf{q}_i is an image \mathbf{x}_i mapped into a query latent via a query encoder $\mathbf{q}_i = f(\mathbf{x}_i)$ The goal is to learn $f(\cdot)$
- \mathbf{k}_i is an image mapped into a key latent via a key encoder $g(\mathbf{x}_i)$ (+ and - are positive and negative samples)
- N is the batch size
- K is the number of negatives per positive sample
- τ is a temperature for the Softmax function



- fMRI data is very high dimensional and many studies have low sample sizes
- In the image domain, Wang et al. 2022³¹ have proposed that salient features should lie on a low-dimensional manifold
- A low rank prior may be useful for our domain

$$\mathcal{L}_i = -\frac{1}{M-1} \sum_m^{M-1} \log \frac{e^{\mathbf{q}_{i,m}^\top \mathbf{k}_i^+ / \tau} \cdot h(\mathbf{Q})}{e^{\mathbf{q}_{i,m}^\top \mathbf{k}_i^+ / \tau} \cdot h(\mathbf{Q}) + \sum_{j=1}^K e^{\mathbf{q}_{i,m}^\top \mathbf{k}_{i,j}^- / \tau}} \quad (5)$$

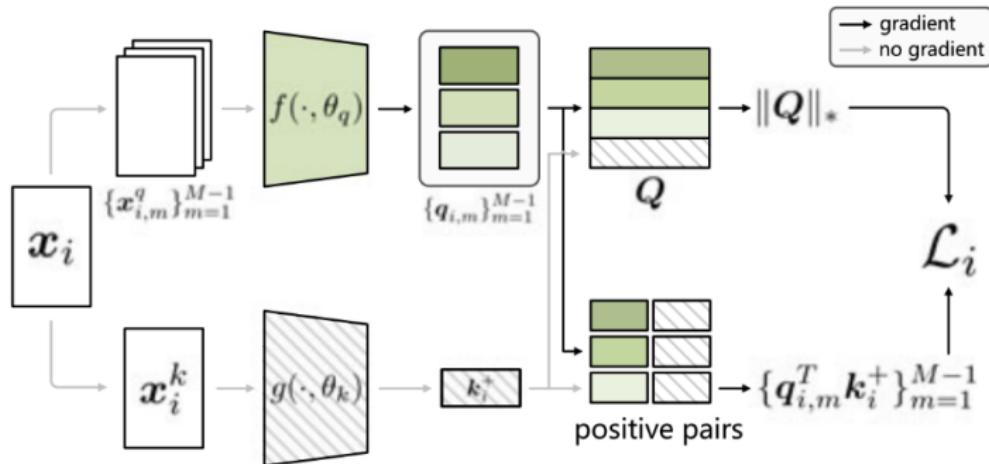
$$h(\mathbf{Q}) = e^{-\frac{\|\mathbf{Q}\|_*}{M \cdot \beta \cdot \tau}}$$

- The Low Rank Promoting Prior for Contrastive Learning (LORAC) modification to the InfoNCE loss



³¹ Wang et al. 2022 10.1109/TPAMI.2022.3180995

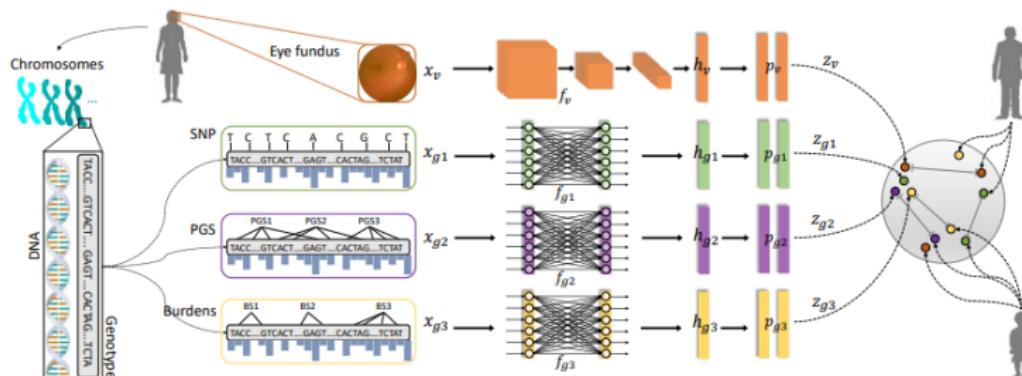
LORAC Training Schematic Diagram



The training for LORAC consists of constructing the matrix \mathbf{Q} out of augmented views of image \mathbf{x}_i and is based on another framework using stop gradients and augmentations

ContIG: Contrastive Learning for Medical Imaging and Genetics

- Contrastive learning has recently been used to integrate medical imaging and genomics³²
- In this study, the medical imaging was of retinal images for diabetic retinopathy



³²Taleb et al. 2022 ContIG CVPR

Contrastive Learning and Multi-Omics Fusion Takeaways

Previous Aims Provide Foundation

- We can leverage work done in previous aims, especially the similarity kernel and augmentations, as a starting point
- Goal is to both
 - ▶ use large amounts of unlabeled data
 - ▶ move multi-omics and brain imaging to a common latent space
- ImageNomer software will be useful for joining and exploring many large or small datasets



Summary and Acknowledgements



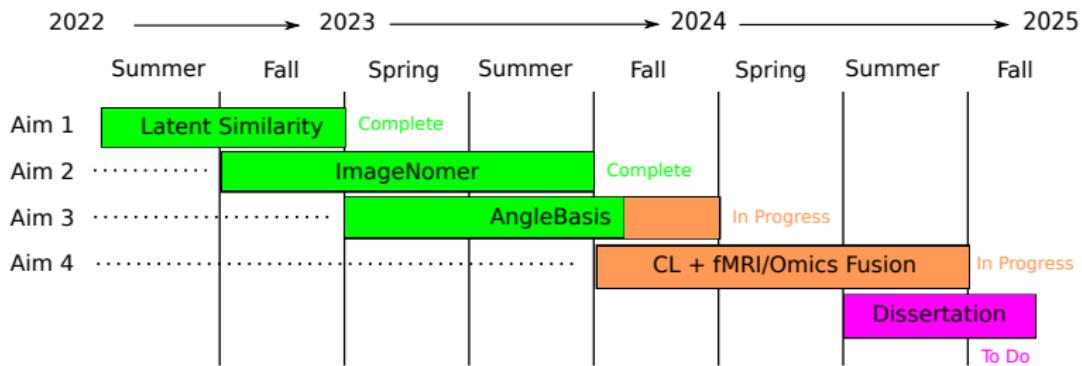
Summary

- The **Latent Similarity** framework works well for small sample size, high dimensionality data
- We apply a greedy feature selection algorithm for better feature selection
- **ImageNomer** allows for quick data exploration, feature visualization, and quality control
- We use ImageNomer to identify an unexpected race confound in fMRI data
- **AngleBasis** provides a generative model for the decomposition and compression of FC
- Our future goal is to leverage **contrastive learning** for fMRI-omics fusion and large unlabeled pre-training



Timeline

- Many of the specific aims have already been started or completed
- In addition, a large amount of data has been acquired from original sources and pre-processed (PNC, BSNIP, ADNI, UKB, etc.)



Acknowledgements



- Lab members: Yu-Ping Wang (PI), Anton Orlichenko, Gang Qu, Binish Patel, Ziyu Zhou (PhD students)
- We would like acknowledge the NIH (grants R01 GM109068, R01 MH104680, R01 MH107354, P20 GM103472, R01 EB020407, R01 EB006841, R56MH124925) and NSF (#1539067) for partial funding support

