

# Task-based fMRI Timeseries Modeling for Alzheimer's Disease Diagnosis

Javier Salazar Cavazos

Department of Electrical Engineering and Computer Science

University of Michigan

April 2<sup>nd</sup>, 2025

# Background on Functional MRI & Alzheimer's Disease

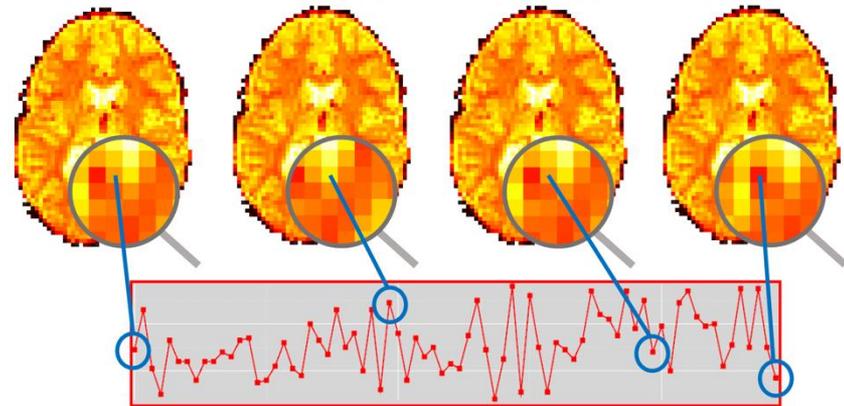
# Background on functional MRI

<https://www.monash.edu/researchinfrastructure/mbi/facilities/human/3t-mri>



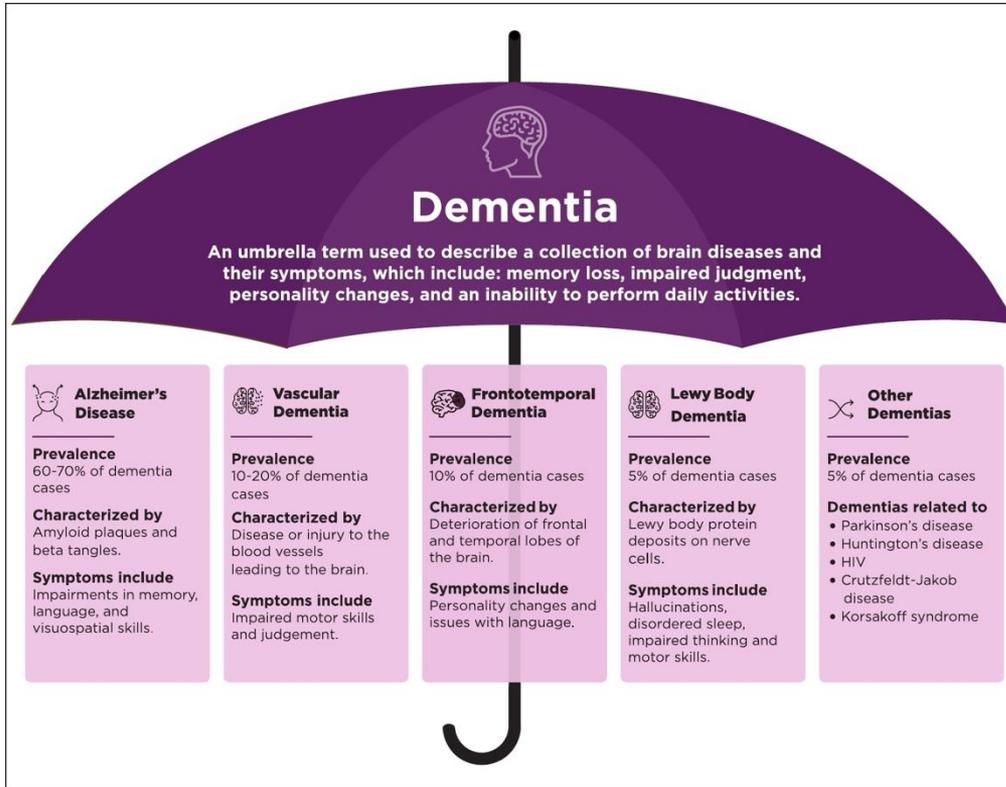
MRI Scanner

[https://martynmcfarquhar.github.io/NCCN-IA-fMRIPreProcessing/1\\_fmri-structure.html](https://martynmcfarquhar.github.io/NCCN-IA-fMRIPreProcessing/1_fmri-structure.html)  
Volume 1 ... Volume 30 ... Volume 60 ... Volume 74



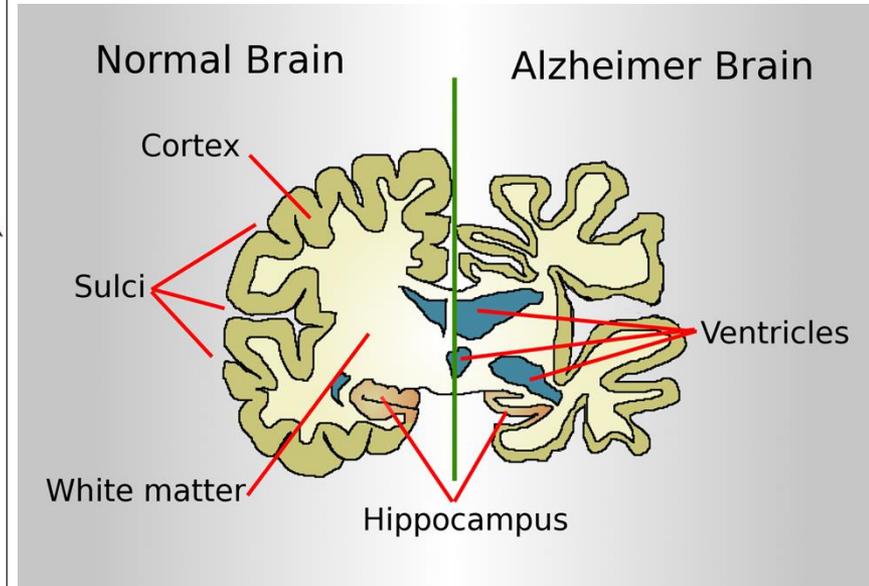
- 3D MRI volume acquired over time to form timeseries for each voxel region
- BOLD data measures blood concentration in gray matter regions
- We use BOLD data as a **proxy** for neural activity

# Background on AD



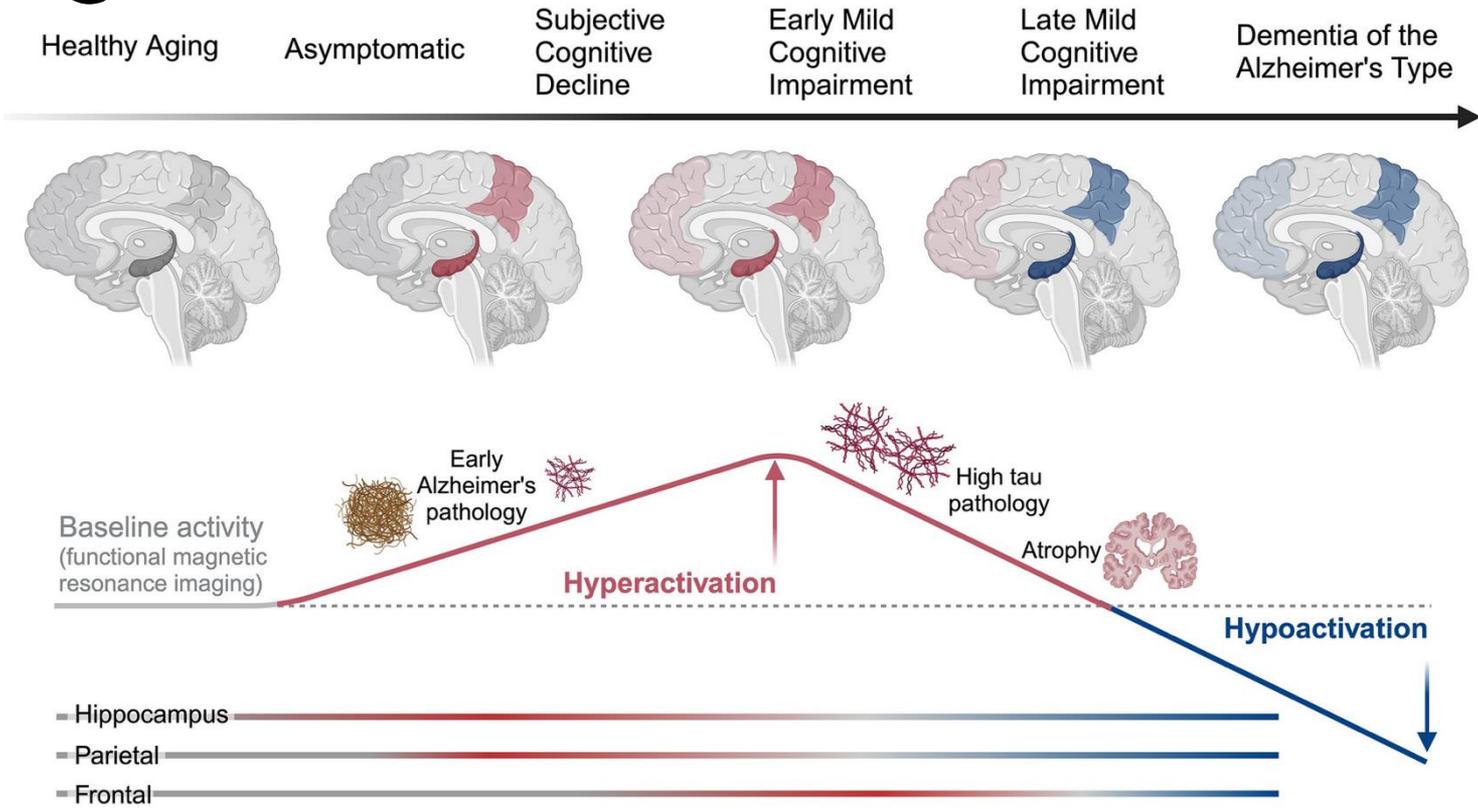
<https://www.summerfieldredlands.com/understanding-the-umbrella-of-dementia/>

## Effects of AD on structural brain



<https://www.emf.ethz.ch/en/knowledge/topics/health/neurodegenerative-diseases/>

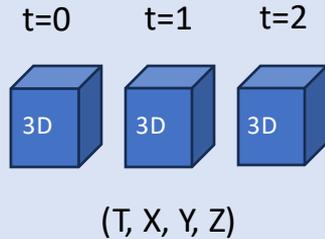
# Background on AD in fMRI



Fischer, Larissa, et al. "Precuneus activity during retrieval is positively associated with amyloid burden in cognitively normal older APOE4 carriers." *Journal of Neuroscience* 45.6 (2025).

# Background - fMRI Data “Views”

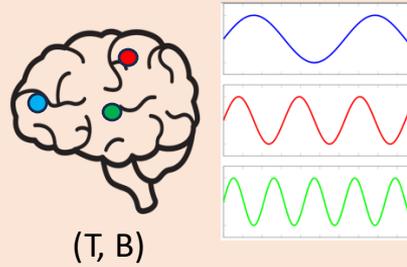
## 4D Spatial-Temporal Data



- High-dimensional
- Difficult to train models
- Harder to draw insights
- No information is lost since it is the raw data

Chp. 5 – 4D CNN w/ rs-fMRI

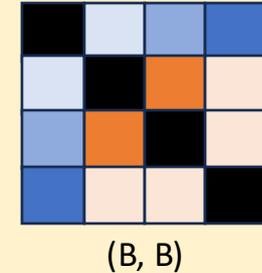
## Timeseries Data



- Low-dimensional
- Highly interpretable
- Easier to train
- Some info lost (especially if  $B$  is small)

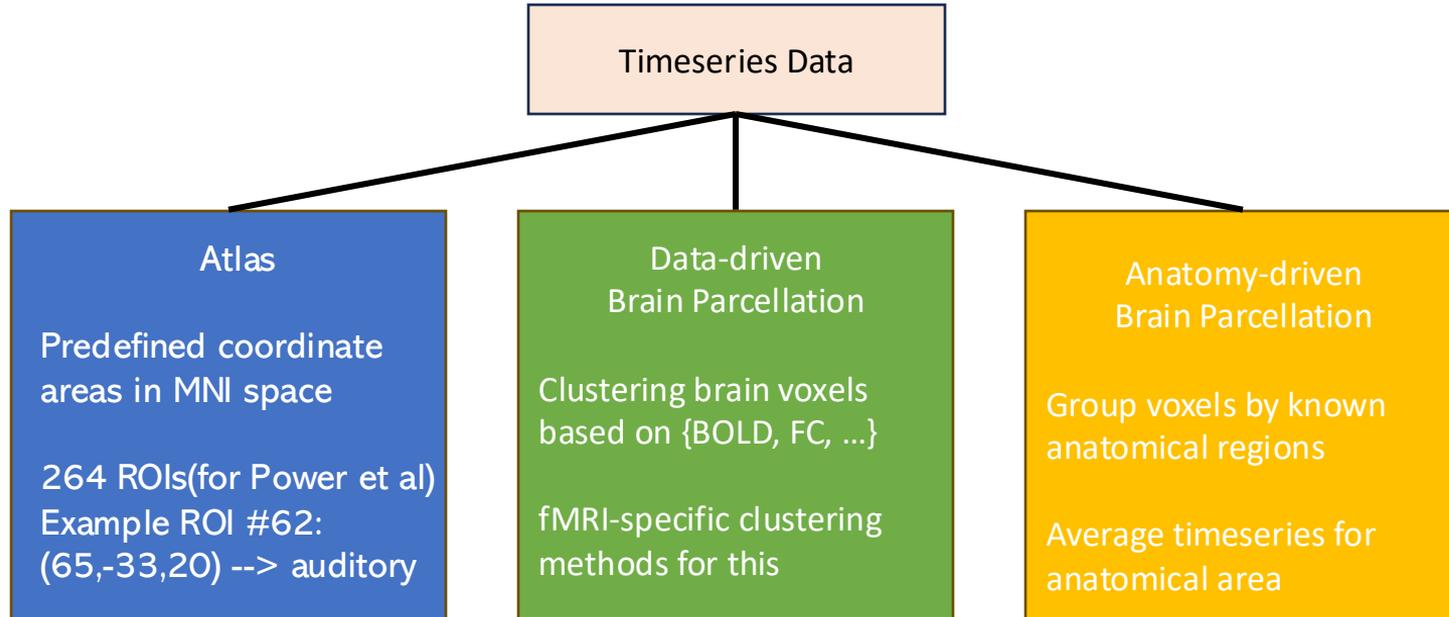
Chp. 8 – Task fMRI

## Functional Connectivity Data



- Super Low-dimensional
- Looks at relationships rather than BOLD activity
- Time information lost entirely

# Timeseries Extraction



For our work, we use a modified “Power” atlas that contains a few extra ROIs (cerebellum)

# **Task-based** fMRI Timeseries Modeling for Alzheimer's Disease **Diagnosis**

# Preliminaries

- Previous work (ISMRM 2025) was on **resting-state** fMRI data that is plentiful due to ADNI consortium
- This ongoing work is on **task-based** fMRI data that is severely more constrained (**zero public datasets for AD**)
- Because of this, the 4D CNN model (ISMRM 2025) is not appropriate and not very interpretable!
- We will require a new approach for this problem

# Preliminaries

- Many meta-analysis papers on MRI/PET/fMRI/... in AD
- Plentiful publications on AD + ML in **resting-state** fMRI
- **Zero** published works on AD + ML in **task-based** fMRI
- There are a few papers that do manual data analysis on task-based fMRI to draw conclusions
- Interesting topic since we believe that task-based fMRI can better “stress” the brain networks relative to resting

# MADC Data (private dataset)

Resting (rs-fMRI)		
CN	MCI	DAT
256	51	41

Object Localization (tb-fMRI)		
CN	MCI	DAT
92	14	10

Face Name Association (tb-fMRI)		
CN	MCI	DAT
183	39	26

Numbers are # of subjects. CN = cognitively normal, MCI = mild cognitive impairment, DAT = Dementia of the Alzheimer's type

- Spatial: 2.4 mm
- Temporal: 0.8s TR, ~5 minute session
- Multi-Band Single-Echo EPI
- 3 scores based on written exams for each subject

For MCI & DAT classes:

We are severely **data constrained** for task-based fMRI

# Goals

- This is not so much about **automated** detection
- Nothing to compare against in **task-based** fMRI
- More about comparing **resting-state** vs **task-based** fMRI as **modalities useful in AD** characterization
- Gaining **insights** as to what ROIs and subnetworks are important for diagnosis in **task-based** fMRI

# Problem

- Explore the potential of task-based fMRI as a biomarker
- Severely **data constrained** in MCI/DAT subjects
- Direct classification models may not be ideal
- “Forced” to use healthy controls only for training
- We will look at alternative problem formulations

# Proposed Ideas

## Pretrain + Finetune

- Use rs-fMRI from ADNI & MADC for classification
- Finetune on task data
- Worth trying due to simplicity

## Data Augmentation

- VAE-style generation
- BioDiffusion; synthesize multivariate class-specific samples
- Not enough samples for BioDiffusion and VAE gives poor resolution
- Does the small training dist. generalize?

## ICA-style Learning

- Learn mixing coefficients to be used for a simple SVM/MLP classifier
- Use or develop a “Deep” ICA (nonlinear mixing) method
- Nonlinear mixing might be better suited for Task-based fMRI data?

## Outlier Detection

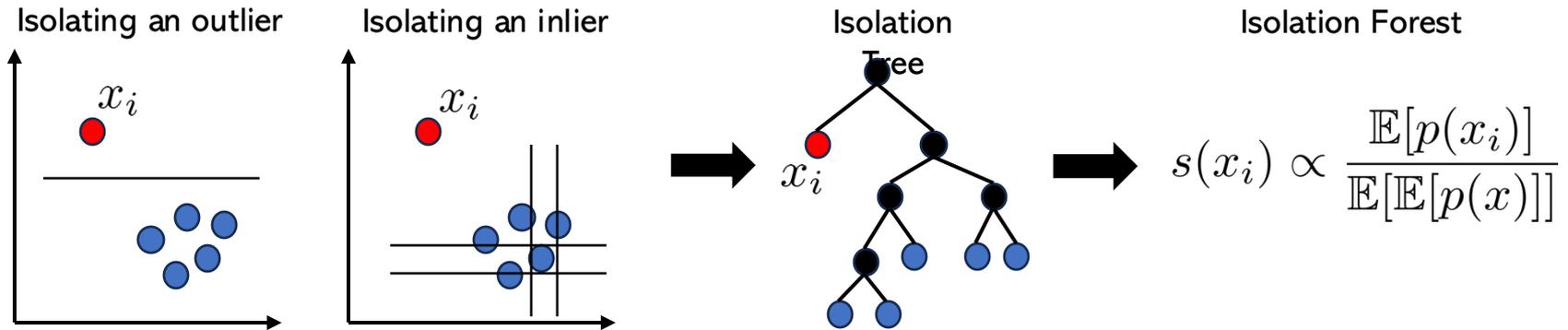
- We have plenty of healthy rs & tb fMRI data
- Train a model to do forecasting/masking/recon/regression
- Use outlier scores to predict severity of disease
- Use behavior score regression since it correlates with disease

We think outlier detection is a good approach so we can use MCI/DAT data for testing only

# Outlier Detection

# Outlier Detection

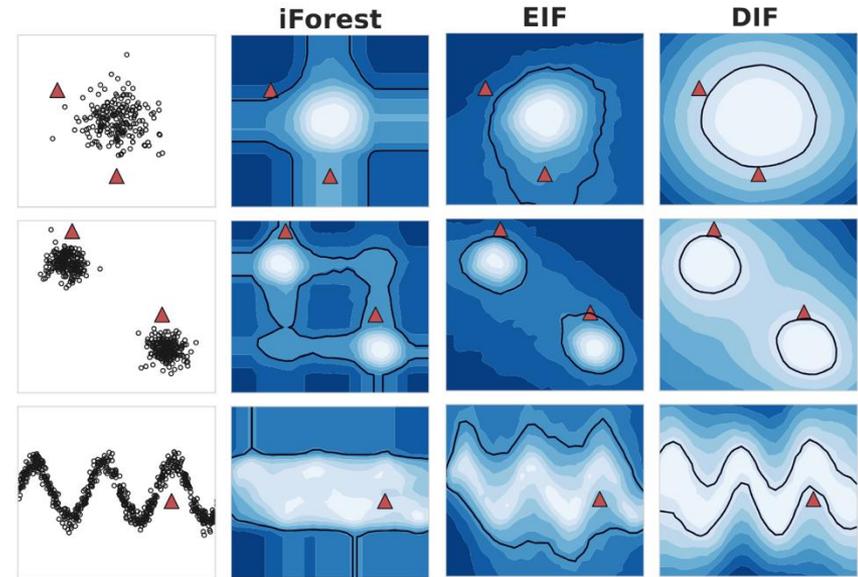
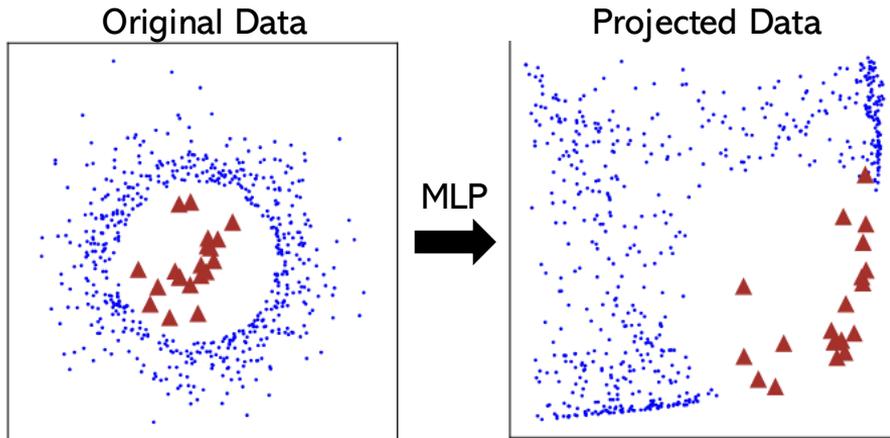
Isolation Forests are one of the most well-known and successful methods used in practice



This is an effective method but limited since it is only vertical/horizontal lines. Is there something better?

# Deep Isolation Forests

- Despite the name, no training!
- Just do isolation forest in random “subspaces”
- Forest = ensemble of trees w/ different projections

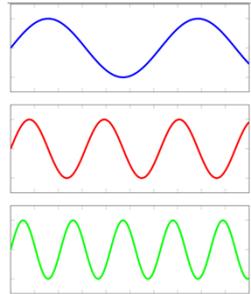
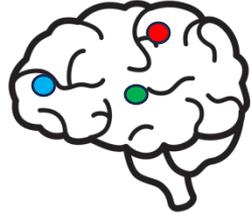


Xu, Hongzuo, et al. "Deep isolation forest for anomaly detection." *IEEE Transactions on Knowledge and Data Engineering* 35.12 (2023): 12591-12604.

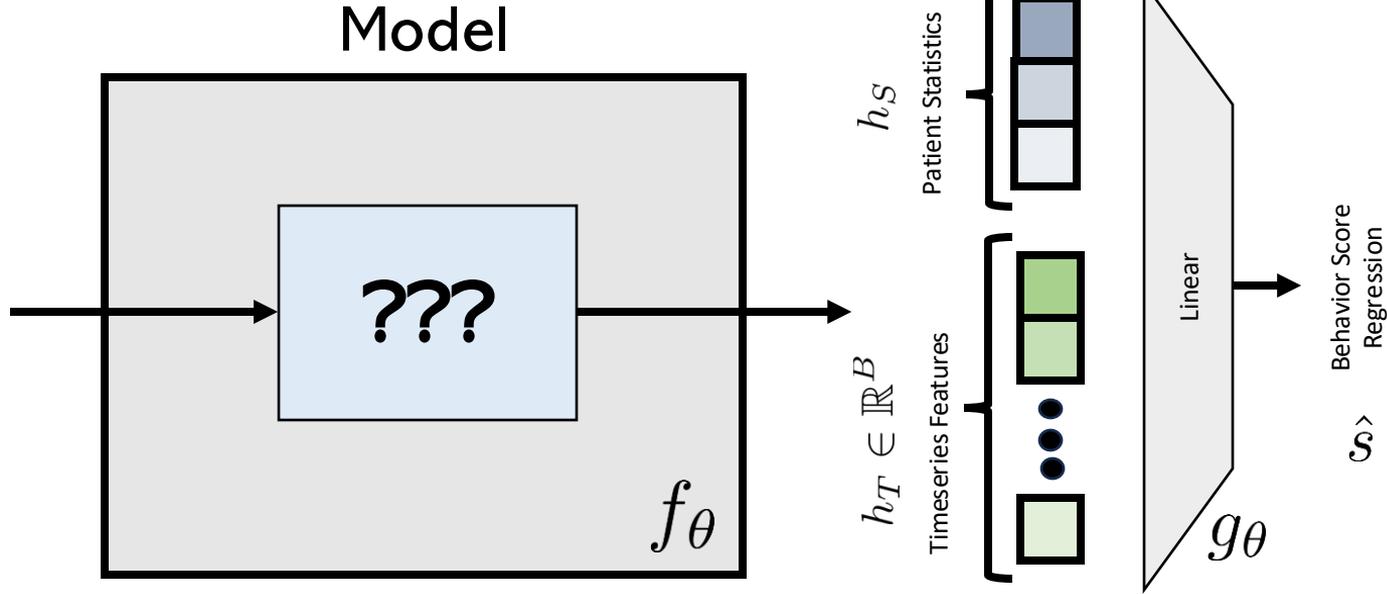
# Outlier Detection

- We hope that our MCI/DAT data is well-separated
- If not, Deep Isolation Forest is a more robust approach
- The question becomes what features to use in DIF?
- The BOLD timeseries data? Many downsides to this
- Better to extract features from some model

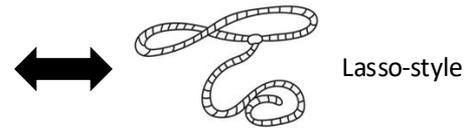
# Score Regression



$$X_T \in \mathbb{R}^{T \times B}$$



$$\min_{\theta} \underbrace{\|s - g_{\theta}(f_{\theta}(X_T), h_s)\|_2^2}_{\text{Data Fidelity}} + \lambda \underbrace{\|f_{\theta}(X_T)\|_1}_{\text{Sparse ROIs}}$$



# Deep State Space Models (SSMs)

# State Space Models (SSMs)

Given input signal  $u(t)$ , generate output sequence  $y(t)$  using state  $x(t)$

$$\left. \begin{aligned} x'(t) &= Ax(t) + Bu(t) \\ y(t) &= Cx(t) + Du(t) \end{aligned} \right\} \text{This is an LTI system}$$

A = State Matrix

B = Input Matrix

C = Output Matrix

D = Feedthrough Matrix

# Toy Example

Consider a mass attached to the wall with a spring.

Let  $m$  = mass,  $k$  = spring constant, and  $b$  = friction constant

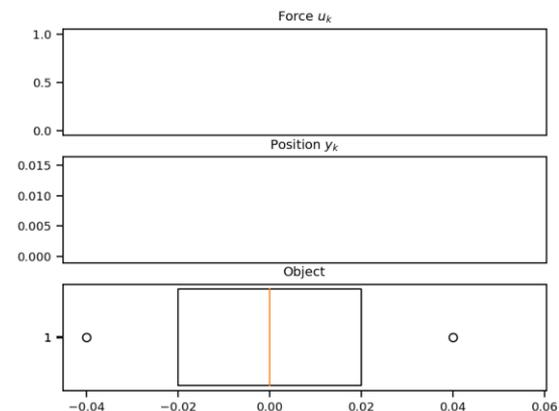
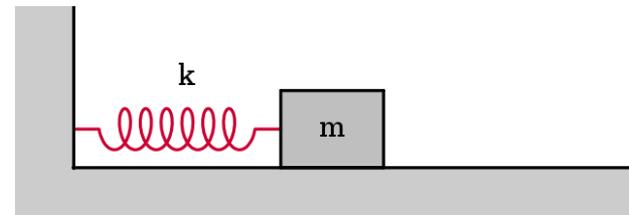
Recall that  $u(t)$  is force applied to spring and  $y(t)$  is position of object

$$my''(t) = u(t) - by'(t) - ky(t)$$



$$A = \begin{bmatrix} 0 & 1 \\ -\frac{k}{m} & -\frac{b}{m} \end{bmatrix}, B = \begin{bmatrix} 0 \\ \frac{1}{m} \end{bmatrix}, C = [1 \quad 0]$$

Of course, in contrast, Deep SSMs learn  $\{A,B,C\}$  for complex systems



# Discretization\*

So far, we have focused on continuous time SSMs, but sampled data is discrete!

Recall differential equations (Euler's method):

$$x'(t) = \lim_{\Delta \rightarrow 0} \frac{x(t + \Delta) - x(t)}{\Delta} \Rightarrow \Delta x'(t) = x(t + \Delta) - x(t)$$

Using this we can say:

$$x(t + \Delta) = \Delta x'(t) + x(t) \quad \text{and recall} \quad x'(t) = Ax(t) + Bu(t)$$

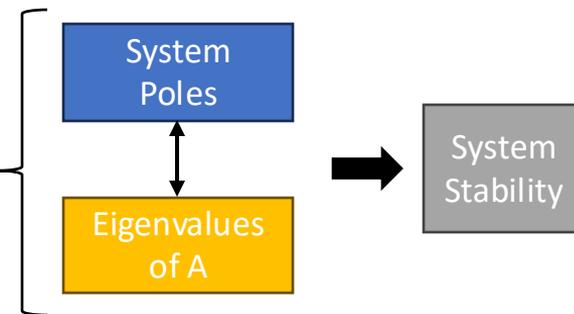
If we combine both from above:

$$x(t + \Delta) = \underbrace{(\Delta A + I)}_{\bar{A}} x(t) + \underbrace{(\Delta B)}_{\bar{B}} u(t)$$

\* Better methods for this exist like zero-order hold (ZOH)

# A: System Matrix

Turns out choice of A/B matrices matter greatly for performance (not surprising!)



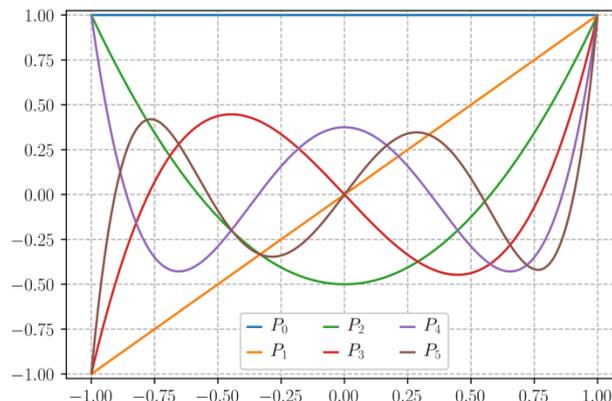
In previous work, it was found that using Legendre polynomials for initialization is a good strategy

$$\underbrace{A_{nk}}_{\text{HiPPO}} = \begin{cases} (2n + 1)^{\frac{1}{2}} (2k + 1)^{\frac{1}{2}} & n > k \\ n + 1 & n = k \\ 0 & n < k \end{cases}$$

Gu, Albert, et al. "How to train your hippo: State space models with generalized orthogonal basis projections." *arXiv preprint arXiv:2206.12037* (2022).

Deep SSM + Structured A = S4

Gu, Albert, Karan Goel, and Christopher Ré. "Efficiently modeling long sequences with structured state spaces." *arXiv preprint arXiv:2111.00396* (2021).



[https://en.wikipedia.org/wiki/Legendre\\_polynomials](https://en.wikipedia.org/wiki/Legendre_polynomials)

# Structured State Space for Sequence Modeling (S4)

Inference:

$$x_k = \bar{A}x_{k-1} + \bar{B}u_k, \quad y_k = \bar{C}x_k$$

Training:

$$y_k = \bar{C}\bar{A}^k\bar{B}u_0 + \bar{C}\bar{A}^{k-1}\bar{B}u_1 + \dots + \bar{C}\bar{B}u_k$$

$$\bar{K} = (\underbrace{\bar{C}\bar{B}, \bar{C}\bar{A}\bar{B}, \dots, \bar{C}\bar{A}^k\bar{B}, \dots}_{\text{1D elements}}) \Rightarrow y = u * \bar{K}$$

Recall that we have an LTI system:

$$y = u * \bar{K} \Leftrightarrow \tilde{Y} = \tilde{U}\tilde{K}$$

We can easily parallelize things for fast training!

In practice,  $\bar{K}$  is not actually constructed. There are some tricks for this.

# Why do SSMs matter?

	sCIFAR
Transformer	62.2
LSTM	63.01
r-LSTM	72.2
UR-LSTM	71.00
UR-GRU	74.4
HiPPO-RNN	61.1
LMU-FFT	-
LipschitzRNN	64.2
TCN	-
TrellisNet	73.42
CKConv	63.74
LSSL	<u>84.65</u>
<b>S4</b>	<b>91.13</b>

Vectorized image data.

	RAW	0.5×
Transformer	<b>X</b>	<b>X</b>
Performer	30.77	30.68
ODE-RNN	<b>X</b>	<b>X</b>
NRDE	16.49	15.12
ExpRNN	11.6	10.8
LipschitzRNN	<b>X</b>	<b>X</b>
CKConv	71.66	<u>65.96</u>
WaveGAN-D	<u>96.25</u>	<b>X</b>
LSSL	<b>X</b>	<b>X</b>
<b>S4</b>	<b>98.32</b>	<b>96.30</b>

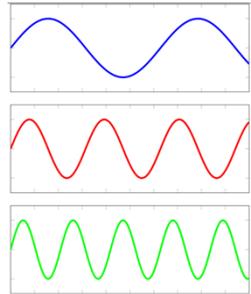
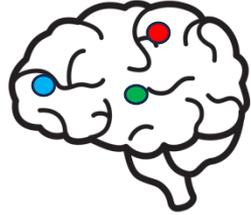
Speech classification.

- Performs well for sequence data
- Vanilla CNNs & Transformers perform poorly with timeseries
- Discretized matrices mean flexibility to sampling rate effects
- Trains like a transformer but does inference like an RNN

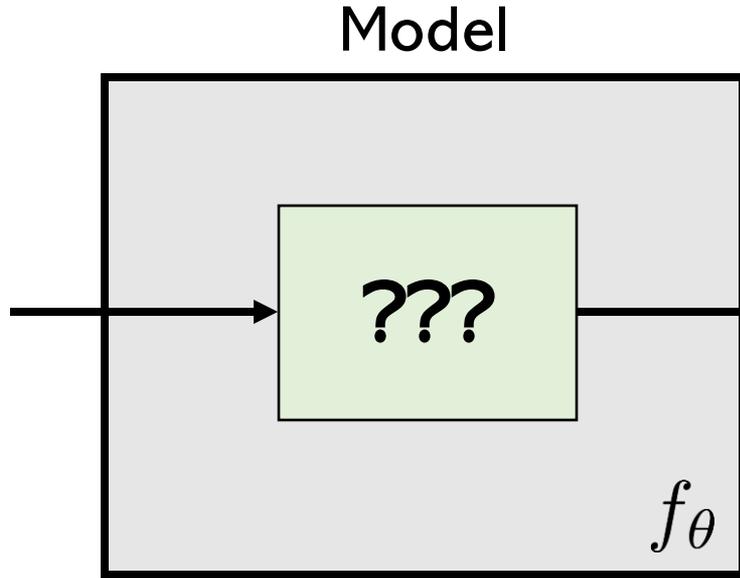


# Score Regression

# Score Regression



$$X_T \in \mathbb{R}^{T \times B}$$

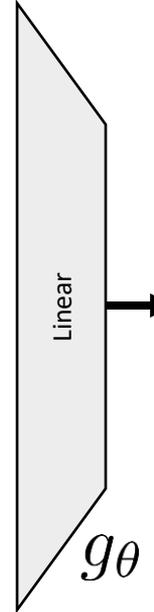


$h_S$

Patient Statistics

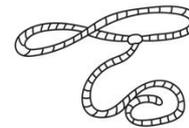
$h_T \in \mathbb{R}^B$

Timeseries Features



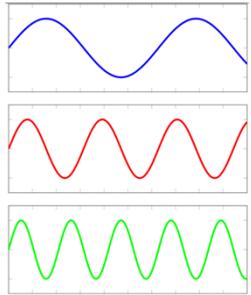
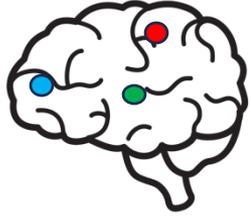
Behavior Score  
Regression

$$\min_{\theta} \underbrace{\|s - g_\theta(f_\theta(X_T), h_s)\|_2^2}_{\text{Data Fidelity}} + \lambda \underbrace{\|f_\theta(X_T)\|_1}_{\text{Sparse ROIs}}$$

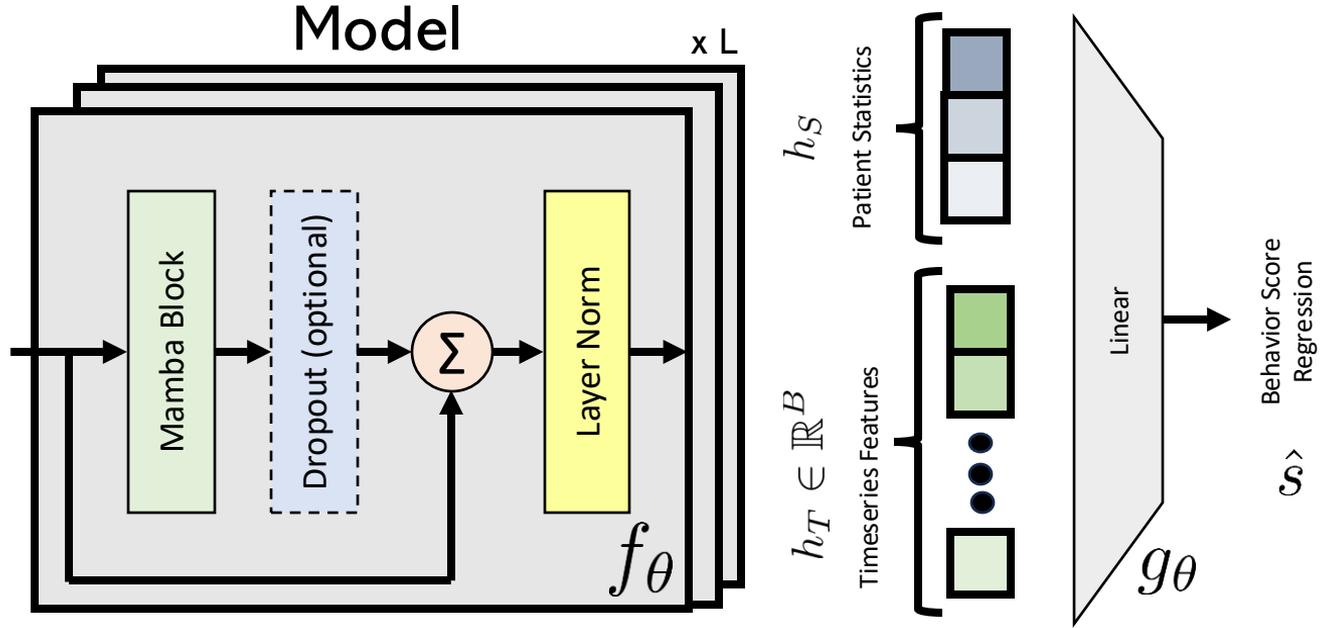


Lasso-style

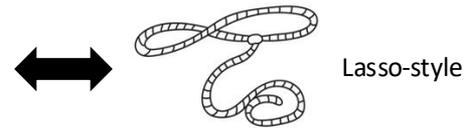
# Score Regression



$$X_T \in \mathbb{R}^{T \times B}$$



$$\min_{\theta} \underbrace{\|s - g_\theta(f_\theta(X_T), h_s)\|_2^2}_{\text{Data Fidelity}} + \lambda \underbrace{\|f_\theta(X_T)\|_1}_{\text{Sparse ROIs}}$$



# The “Limitation” of Mamba

Univariate Data

$$x'(t) = Ax(t) + Bu(t)$$

$$y(t) = Cx(t) + Du(t)$$



Multivariate Data

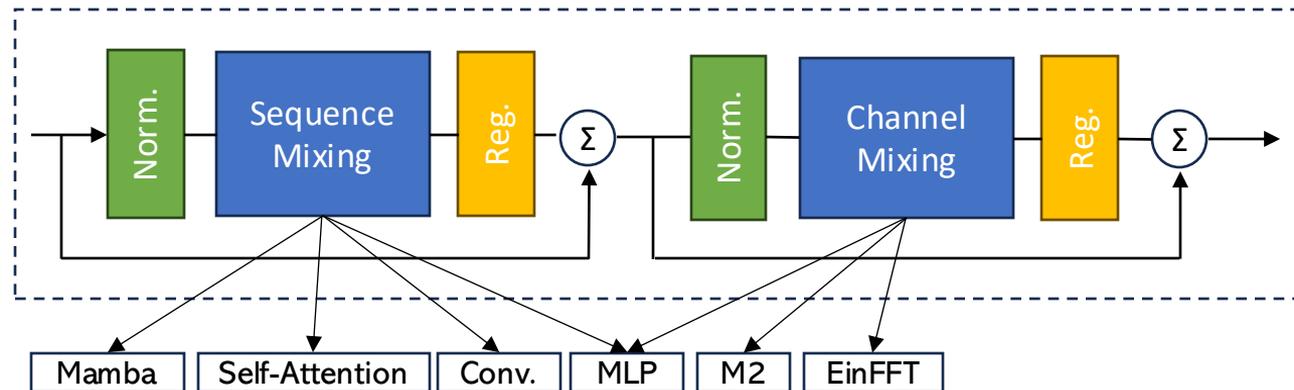
$$\mathbf{x}'(t) = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t)$$

$$\mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t)$$

Linearity; Separate SSMs for each feature dim.

We are capturing intra-variable features (temporal patterns) but not inter-variable features (interactions)!

Modern Model Backbone

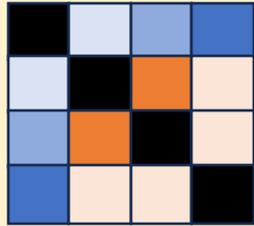


- Transformer = Self-Attention + MLP
- CNN = Convolution + MLP
- MLP-Mixer = MLP + MLP
- Mamba Model = Mamba + MLP

We can use this to learn ROI interactions but at the huge expense of interpretability!

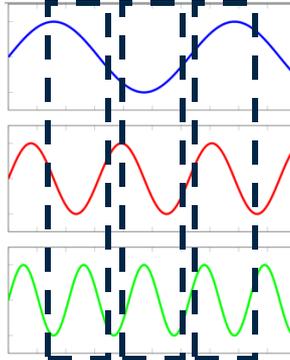
# Learning spatial interactions

## Functional Connectivity



(B, B)

- Super Low-dimensional
- Looks at relationships rather than BOLD activity
- Time information lost entirely

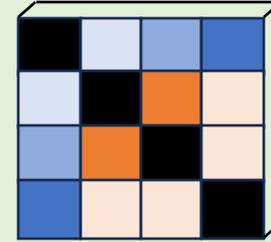


Note:

$B = 272$  so a vectorized graph is quite large!

$$0.5 * B * (B - 1) = \sim 37k$$

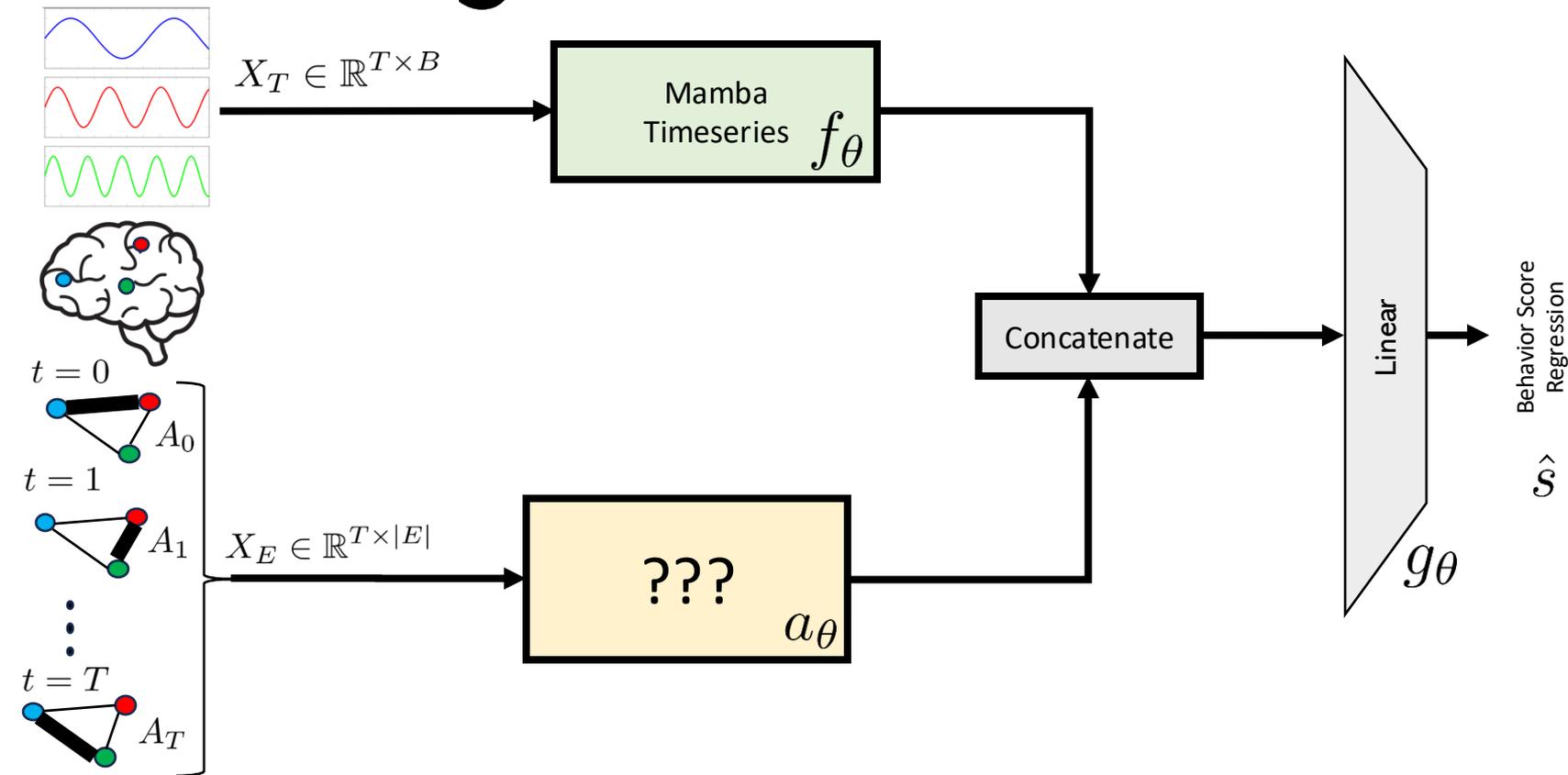
## Dynamic Functional Connectivity



(T, B, B)

- Looks at relationships rather than BOLD activity
- Captures ROI interaction dynamics over time
- Useful for task-based fMRI

# Score Regression



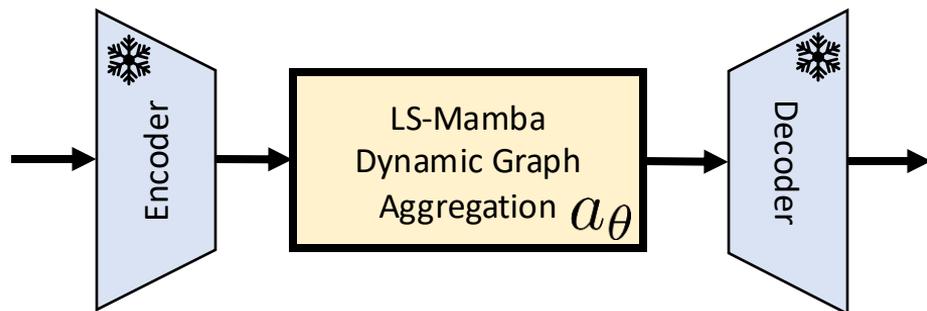
# Dynamic Graph Aggregation

Lets “summarize” all ROI interactions to one useful FC for score regression:

$$X_E \in \mathbb{R}^{T \times |E|} \longrightarrow a_\theta(X_E) \in \mathbb{R}^{|E|}$$

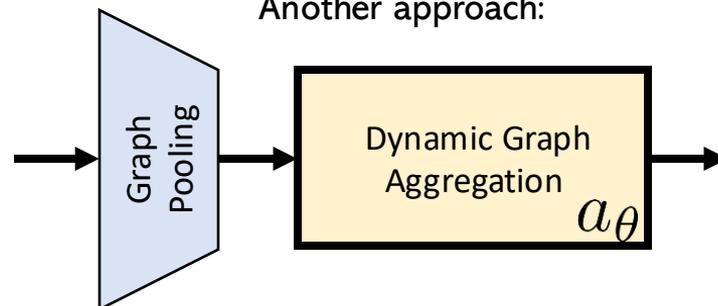
Remember that  $|E| = 37k$  which is very large to model directly!

One possible approach:



- Learn ROI interactions in a compressed space
- Good assumption since nearby ROIs likely have similar activity, i.e., fMRI data is correlated

Another approach:

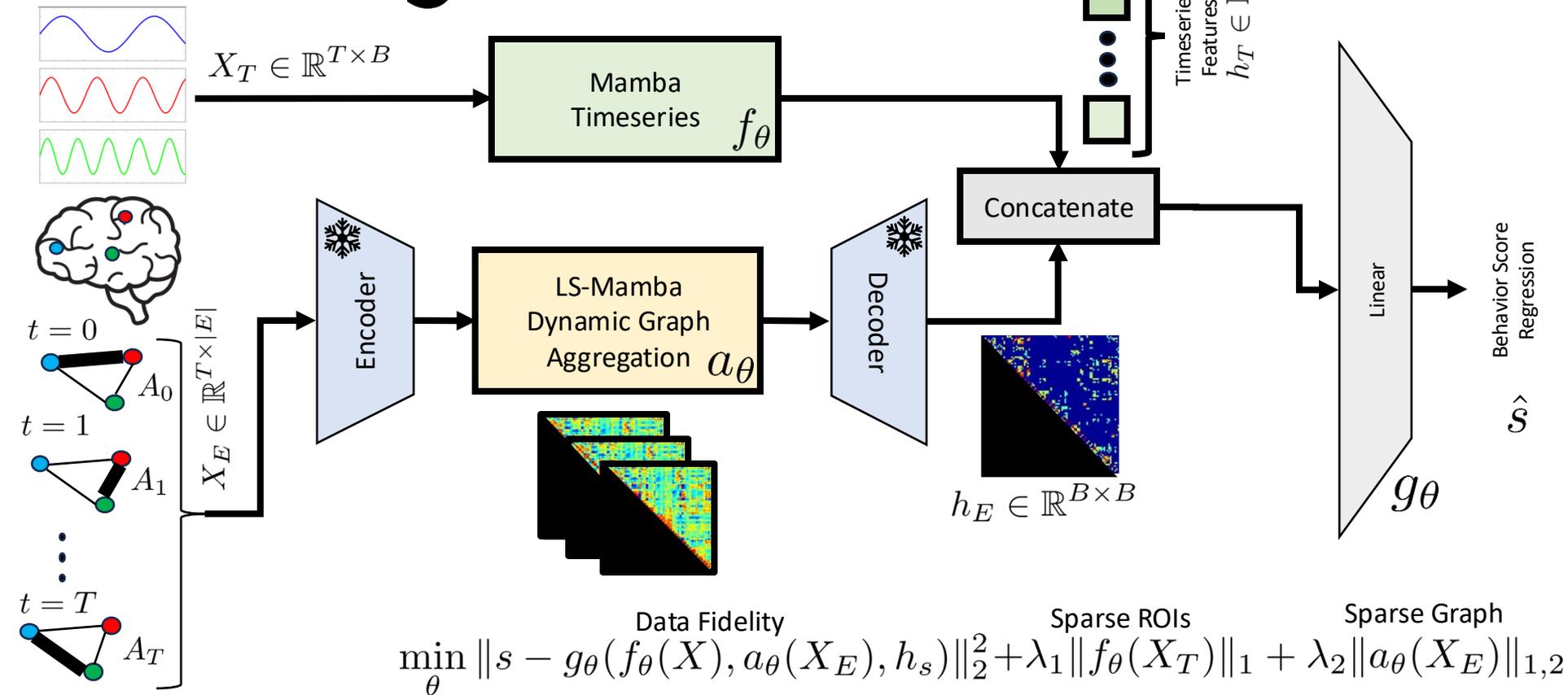


- Top-k pooling can learn to pick the k most important nodes (e.g., drop visual cortex)
- DiffPool can cluster similar nodes together (e.g., merge all ROIs in parietal lobe)

Cangea, Cătălina, et al. "Towards sparse hierarchical graph classifiers." *arXiv preprint arXiv:1811.01287*(2018).

Ying, Zhitao, et al. "Hierarchical graph representation learning with differentiable pooling." *Advances in neural information processing systems* 31 (2018).

# Score Regression



# Cost Function

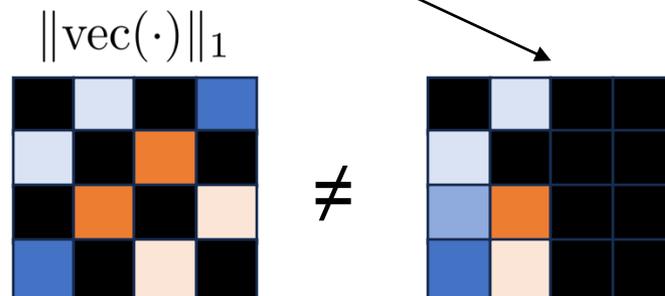
$$\min_{\theta} \overset{\text{Data Fidelity}}{\|s - g_{\theta}(f_{\theta}(X), a_{\theta}(X_E), h_s)\|_2^2} + \overset{\text{Sparse ROIs}}{\lambda_1 \|f_{\theta}(X_T)\|_1} + \overset{\text{Sparse Graph}}{\lambda_2 \|a_{\theta}(X_E)\|_{1,2}}$$

Sparse ROIs  
 $\lambda_1 \|f_{\theta}(X_T)\|_1$



← Visual cortex temporal features probably not useful

Sparse Graph  
 $\lambda_2 \|a_{\theta}(X_E)\|_{1,2}$



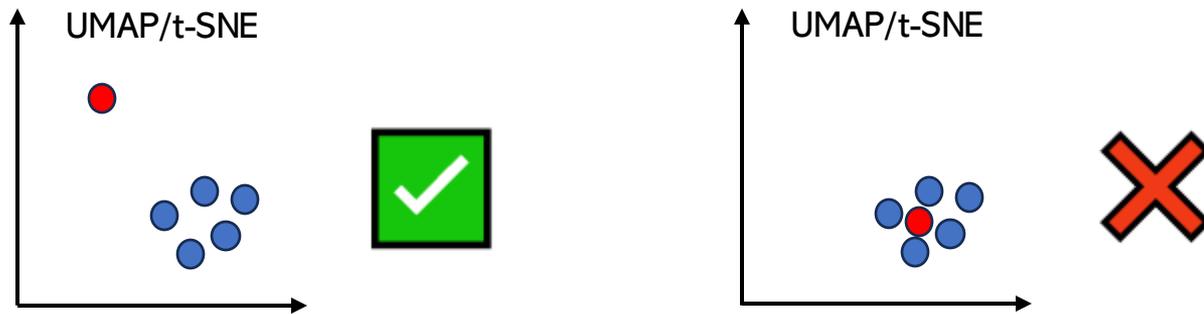
We are removing useless nodes in the graph as opposed to having sparse edges to find important fMRI subnetworks!

# Post-training

Use extracted features for any of the following tasks:

- **Outlier detection (Deep Isolation Forest)**
- Clustering (K-means, PET-TURTLE (mine), etc...)
- Classification (MLP, SVM, K-nearest neighbor)

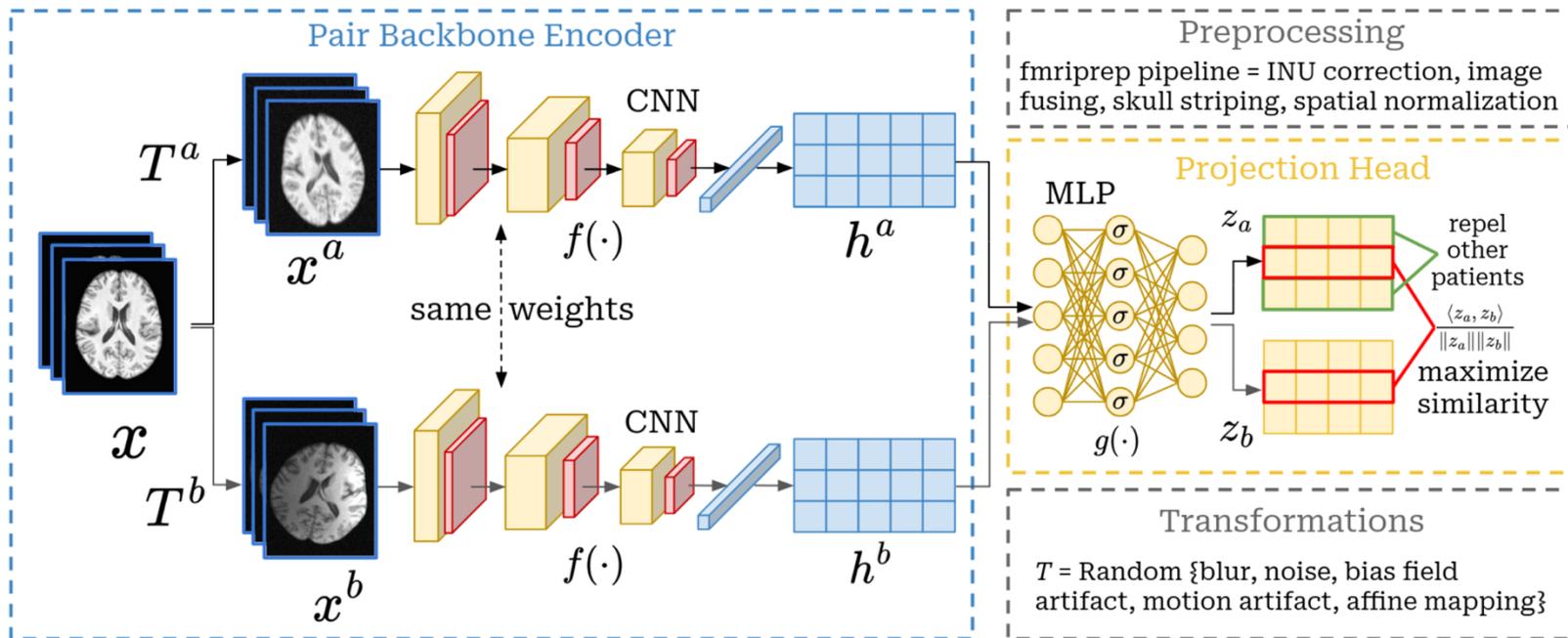
But firstly, we should visualize the extracted features to ensure these approaches are viable:



## Possible Alternative to Score Regression (if time allows)

# Contrastive Learning

Structural MRI ISMRM submission involved contrastive learning  
Could we do something similar in fMRI for unsupervised learning?  
We can explore gender prediction, score regression, etc...



# Contrastive Learning

Method	Accuracy	Sensitivity	Specificity
2 class setting (CN/DAT) ref. classification <b>0.86 ACC</b> (3D CNN)			
CL + LP	<b>0.82</b>	<b>0.91</b>	<b>0.73</b>
VAE + MLP	0.76	0.88	0.62
3 class setting (CN/MCI/DAT) ref. classification <b>0.58 ACC</b> (3D CNN)			
CL + LP	<b>0.56</b>	<b>0.57</b>	<b>0.61</b>
VAE + MLP	0.47	0.48	0.58

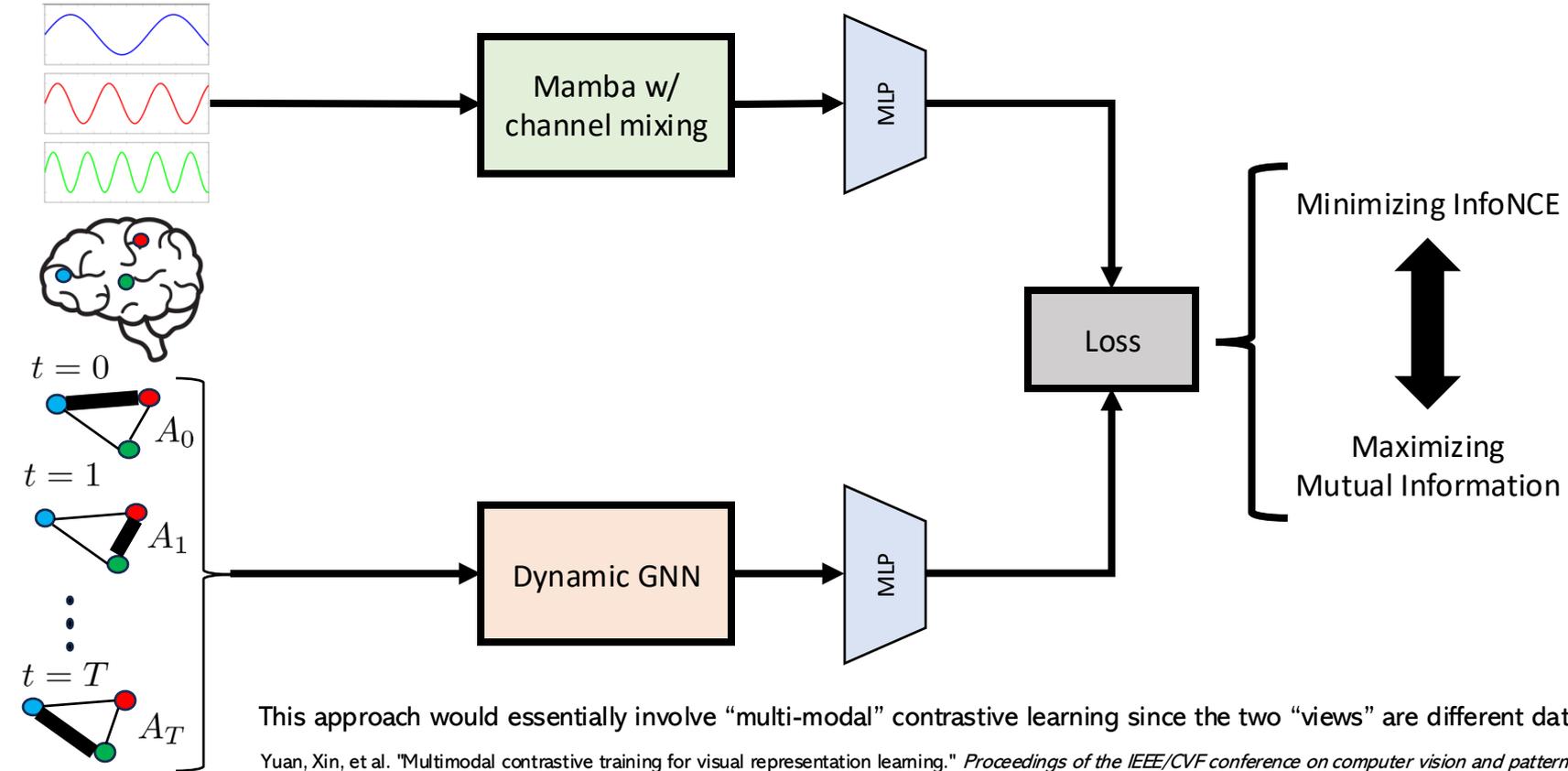
From my sMRI experiments, I found that CL is very competitive (close to end-to-end classification model) relative to VAE approach

Other works such as deep clustering methods using contrastive loss have found the same:

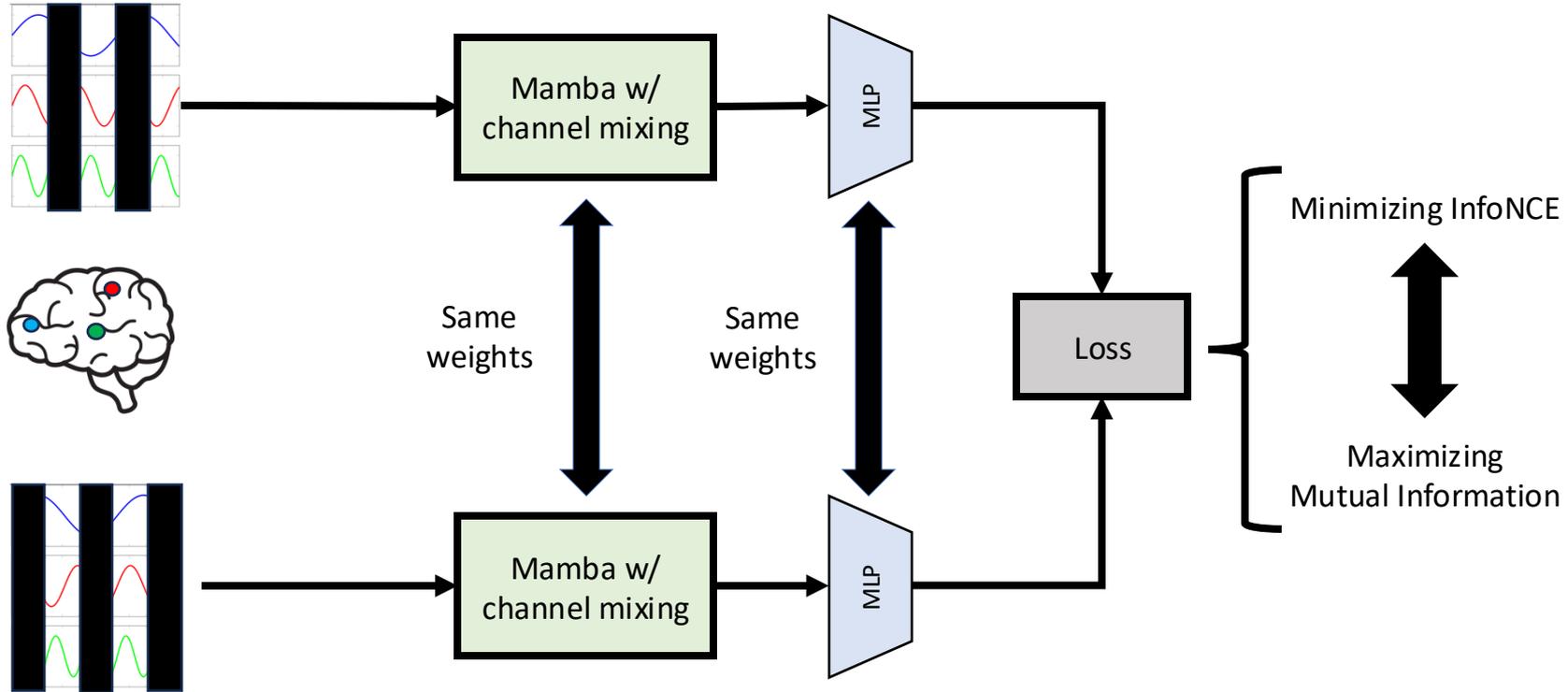
Dataset (ACC)	VAE	DAC = VAE + K-means Loss	Contrastive Clustering (CC)
CIFAR-10	0.291	0.522	<b>0.790</b>
CIFAR-100	0.152	0.238	<b>0.429</b>
STL-10	0.282	0.470	<b>0.850</b>
ImageNet-10	0.334	0.527	<b>0.893</b>

Li, Yunfan, et al. "Contrastive clustering." *Proceedings of the AAAI conference on artificial intelligence*. Vol. 35. No. 10. 2021.

# Possible idea #1



# Possible idea #2



This approach would involve timeseries contrastive learning with some mechanism to create “views” such as masking:

Pöppelbaum, Johannes, Gavneet Singh Chadha, and Andreas Schwung. "Contrastive learning based self-supervised time-series analysis." *Applied Soft Computing* 117 (2022): 108397.

# Conclusion

- **Task-based** fMRI data could be a promising modality
- **Limited data** size --> alternative problem statements
- Score Regression + Outlier Detection
- We separately model ROI timeseries + spatial interactions in an interpretable way
- Possible to infer biological insights about what subnetworks are important for AD

# Thank you for your attention!

Javier Salazar Cavazos

Advisors: Jeffrey A. Fessler, Laura Balzano, and Scott Peltier

Department of Electrical Engineering and Computer Science

University of Michigan

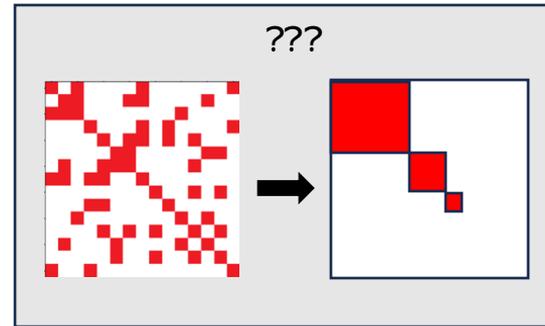
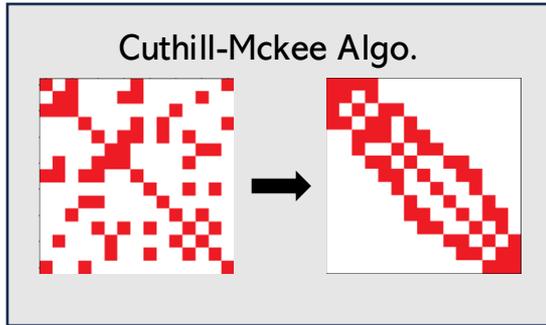
April 2<sup>nd</sup>, 2025

# My Questions

# Q1: Block Diagonal Reordering

$$\min_{\theta} \underbrace{\|s - g_{\theta}(f_{\theta}(X)), a_{\theta}(X_E), h_s\|_2^2}_{\text{Data Fidelity}} + \lambda_1 \underbrace{\|f_{\theta}(X_T)\|_1}_{\text{Sparse ROIs}} + \lambda_2 \underbrace{\|a_{\theta}(X_E)\|_{1,2}}_{\text{Sparse Graph}}$$

This matrix will not be structured!

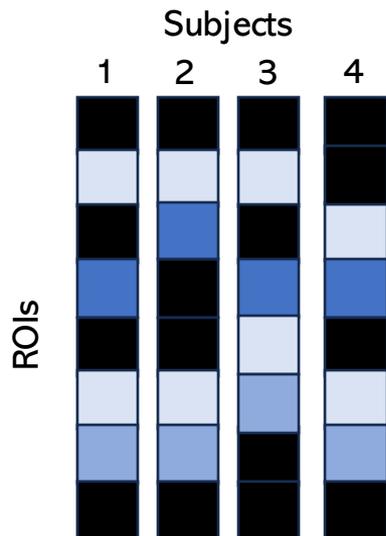


My thoughts:

- Convert correlation matrix to affinity matrix ( [-1, 1] -> [0, bound] , critical to think about function! )
- Apply spectral clustering to find minimum “cuts” in graph
- Use labels to reorder symmetric matrix
- Optionally, look at the eigenspectrum of the graph Laplacian to find # of clusters (unknown here!)

# Q2: Population Variance Control

Prototype enforcement?



How to ensure each subject has same ROIs selected?

My thoughts:

$$k = \text{softmax}(h_T \odot p)$$

$$i = k_i > \lambda \in [0, 1]$$

$$\tilde{h}_T = (h_T \odot k)_i$$

Learnable projection!

Learnable prototype that works across the healthy population  
(removes L1 norm in cost function)