

Interpretable Graph Neural Networks for Connectome-Based Brain Disorder Analysis

Hejie Cui, Wei Dai, Yanqiao Zhu, Xiaoxiao Li, Lifang He, and Carl Yang
(corresponding: j.carlyang@emory.edu)

Department of Computer Science, Emory University, Atlanta, GA 30322, USA



EMORY
UNIVERSITY

CONTRIBUTION HIGHLIGHTS

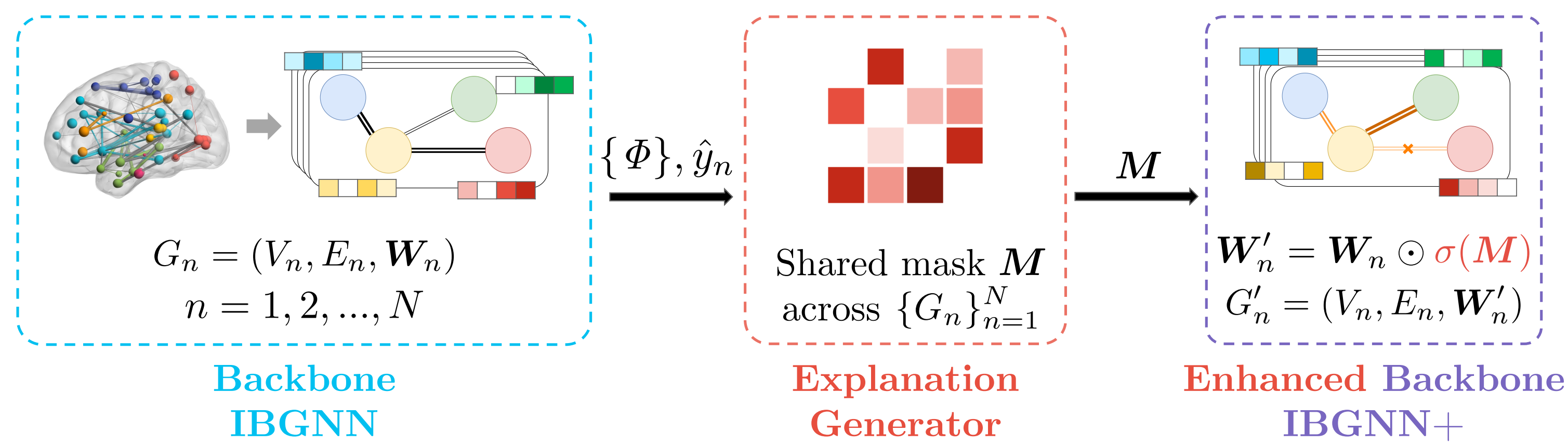


Figure: An overview of our proposed framework.

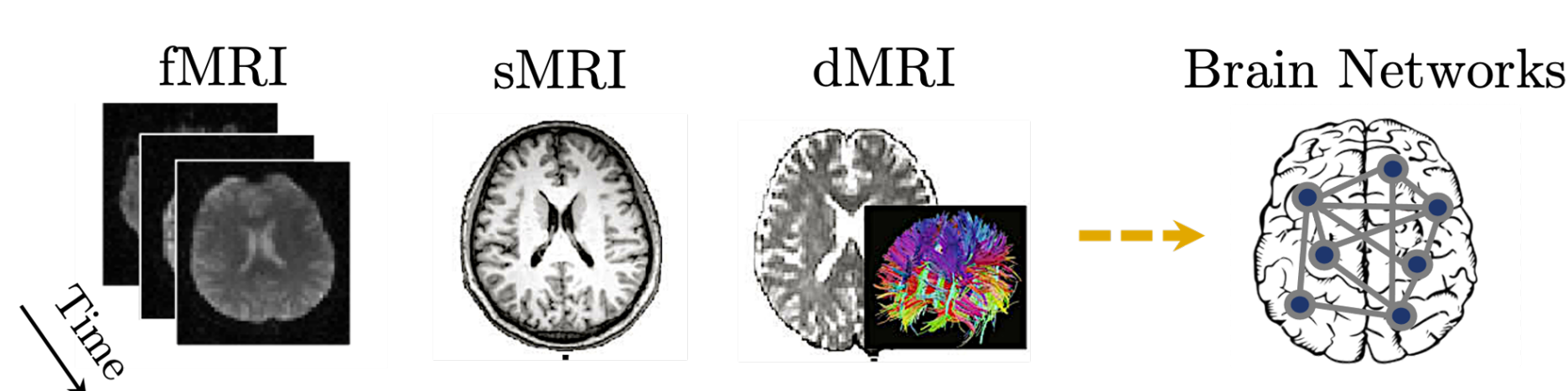
The whole training pipeline of IBGNN+:

- The backbone model is first trained on *the original data*
- Then, the explanation generator learns a *globally shared mask* across subjects
- Finally, we enhance the backbone by *applying the learned explanation mask* and fine-tune the model

INTRODUCTION

Brain Networks

- Brains lie at the core of neurobiological systems
- Mapping the connections of the brain as a network is one of the most pervasive paradigms in neuroscience (Nodes: anatomic regions; Edges: connectivities between the regions)
- Interpretable models on brain networks are vital



Graph Neural Networks (GNNs)

- GNNs have emerged and proved its power for analyzing graph-structured data.
- Compared with shallow models \rightarrow universal expressiveness to capture the sophisticated connectome structures
- However, as a family of deep models, it is prone to **overfitting** and **lack of transparency** and interpretation in predictions!

GNN Explanation

- Existing work mostly focus on general graphs and node-level prediction task and produce a unique explanation for each subject when applied to graph-level tasks
- For brain networks, subjects with the same disorder share similar connection patterns and brain networks possess unique properties
- **Our Motivations:** (1) Unleash the prediction power of GNNs for brain network analysis; (2) Investigate disease-specific patterns common across the group and provide interpretations of different levels

PROBLEM DEFINITION

- **Input:** a set of N weighted brain networks, for each network $G = (V, E, W)$, $V = \{v_{ij}\}_{i=1}^M$ is the Regions Of Interest (ROIs) node set of size M ; $E = V \times V$ is the edge set of brain connectome; $W \in \mathbb{R}^{M \times M}$ is the weighted adjacency matrix describing the connection strengths between ROIs
- **Output:** A prediction \hat{y}_n for each subject n ; A disorder-specific interpretation matrix $M \in \mathbb{R}^{M \times M}$ shared across all subjects, highlighting disorder-specific biomarkers

IBGNN+

Module 1: The Backbone Model IBGNN

- Message Vector: concatenate embeddings of a node v_i , its neighbor v_j , and edge weight w_{ij}
 $m_{ij}^{(l)} = \text{MLP}_1 \left([h_i^{(l)}; h_j^{(l)}; w_{ij}] \right)$
- Propagation Rule
 $h_i^{(l)} = \xi \left(\sum_{v_j \in \mathcal{N}_i \cup \{v_i\}} m_{ij}^{(l-1)} \right)$
- Readout Function: summarize all node embeddings to a graph-level one, with MLP and residual connections
 $z = \sum_{i \in V} h_i^{(L)}, \quad g = \text{MLP}_2(z) + z$
- Training Objectives: cross-entropy
 $\mathcal{L}_{\text{CLF}} = -\frac{1}{N} \sum_{n=1}^N (y_n \log(\hat{y}_n) + (1 - y_n) \log(1 - \hat{y}_n))$

Module 2: The Globally Shared Explanation Generator

- Maximize the agreement between the predictions \hat{y} on the original graph G and \hat{y}' on an explanation graph $G' = (V, E, W')$ induced by a masking matrix M , where $W' = W \odot \sigma(M)$,
 $\mathcal{L}_{\text{MASK}} = -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \mathbb{1}[\hat{y}_n = c] \log P_\Phi(\hat{y}'_n = \hat{y}_n | G'_n)$
- Two regularization terms: encourage the compactness of the explanation and the discreteness of the mask values
 $\mathcal{L}_{\text{SPS}} = \sum_{ij} M_{ij}, \quad \mathcal{L}_{\text{ENT}} = -(M \log(M) + (1-M) \log(1-M))$
- Training Objectives
 $\mathcal{L} = \mathcal{L}_{\text{CLF}} + \alpha \mathcal{L}_{\text{MASK}} + \beta \mathcal{L}_{\text{SPS}} + \gamma \mathcal{L}_{\text{ENT}}$

Enhancing the Backbone with the Learned Explanation: IBGNN+

- Apply the shared global explanation mask to individual brain networks \rightarrow *predictions* and *interpretations* are produced in a closed-loop for brain network analysis

PREDICTION PERFORMANCE

Method	HIV			BP			PPMI		
	Accuracy	F1	AUC	Accuracy	F1	AUC	Accuracy	F1	AUC
M2F	57.14 _{±19.17}	53.71 _{±19.80}	57.50 _{±18.71}	52.56 _{±13.96}	51.65 _{±13.38}	52.42 _{±13.85}	78.69 _{±1.78}	45.81 _{±4.17}	50.39 _{±2.39}
MIC	54.29 _{±8.95}	53.63 _{±19.44}	55.42 _{±19.10}	62.67 _{±20.92}	63.00 _{±21.61}	61.79 _{±21.74}	79.11 _{±2.16}	49.65 _{±3.16}	52.39 _{±2.94}
MPCA	67.14 _{±20.25}	64.28 _{±21.47}	69.17 _{±20.17}	52.56 _{±13.12}	50.43 _{±14.99}	52.42 _{±13.69}	79.15 _{±0.97}	44.18 _{±0.18}	50.00 _{±0.00}
MK-SVM	65.71 _{±7.00}	62.08 _{±7.49}	65.83 _{±7.41}	57.00 _{±8.89}	41.08 _{±13.44}	53.75 _{±10.00}	79.15 _{±0.97}	44.18 _{±0.18}	50.00 _{±0.00}
GCN	70.00 _{±12.51}	68.35 _{±13.28}	73.58 _{±9.49}	55.56 _{±13.86}	50.71 _{±11.75}	61.55 _{±28.77}	78.55 _{±1.38}	47.87 _{±4.40}	59.43 _{±4.44}
GAT	71.43 _{±11.66}	69.79 _{±10.83}	77.17 _{±9.42}	63.34 _{±9.15}	60.42 _{±7.56}	67.07 _{±3.98}	79.02 _{±1.23}	45.85 _{±3.16}	64.40 _{±3.87}
PNA	57.14 _{±12.78}	45.09 _{±19.62}	57.14 _{±12.78}	63.71 _{±11.34}	55.54 _{±14.06}	60.30 _{±11.89}	79.36 _{±1.84}	51.76 _{±10.32}	54.71 _{±6.77}
BrainNetCNN	69.24 _{±10.04}	67.08 _{±11.11}	72.09 _{±10.01}	65.83 _{±20.04}	64.74 _{±17.42}	64.32 _{±13.72}	55.20 _{±12.63}	55.45 _{±9.32}	52.54 _{±10.21}
BrainGNN	74.29 _{±12.10}	73.49 _{±10.75}	75.00 _{±10.56}	68.00 _{±12.45}	62.33 _{±13.01}	74.20 _{±12.93}	69.17 _{±0.00}	44.19 _{±0.00}	45.26 _{±0.65}
IBGNN	82.14 _{±10.81}	82.02 _{±10.36}	86.86 _{±11.65}	73.19 _{±12.20}	72.87 _{±12.29}	83.64 _{±9.81}	79.82 _{±2.47}	51.58 _{±4.66}	70.65 _{±5.35}
IBGNN+	84.29 _{±12.94}	83.86 _{±13.42}	88.57 _{±10.89}	76.33 _{±13.08}	76.13 _{±13.01}	84.61 _{±9.08}	79.55 _{±1.07}	56.58 _{±7.43}	72.76 _{±5.73}

- Backbone IBGNN outperforms shallow/deep baselines (up to 11% absolute improvement)
- The explanation enhanced IBGNN+ further improve the backbone by 9.7% relatively
- IBGNN+ can effectively highlight the disorder-specific signals while achieving the benefit of restraining random noises

INTERPRETATION ANALYSIS

Neural System Mapping

- ROIs on brain networks can be partitioned into different neural systems

I. Salient ROIs

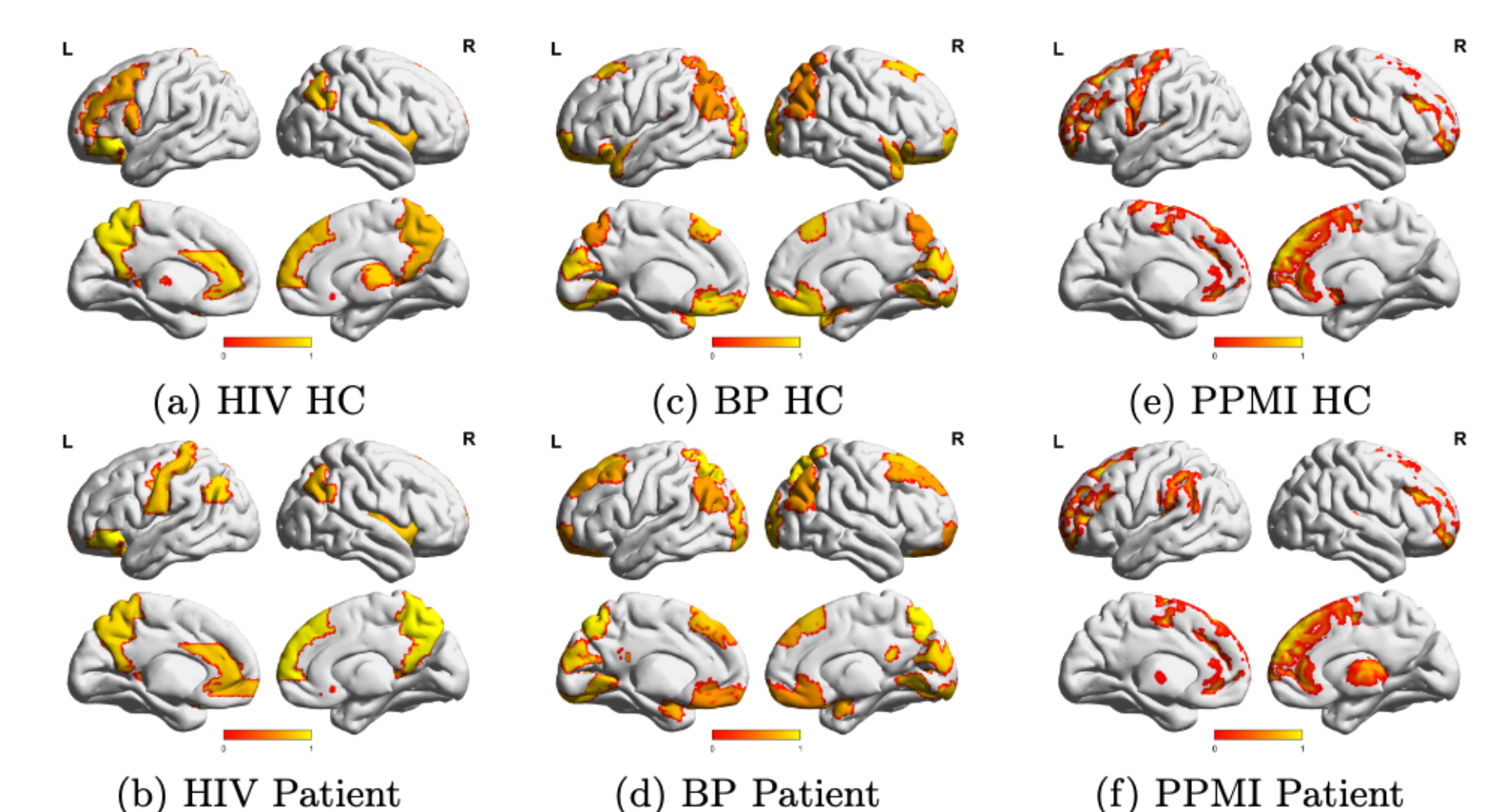


Figure: Salient ROIs on the explanation enhanced brain networks for Health Control (HC) and Patient.

- **Group-level & Individual-level** interpretations on which ROIs contribute most to the prediction of a specific disorder:
HIV: anterior cingulate, paracingulate gyri, inferior frontal gyrus
BP: secondary visual cortex and medial to superior temporal gyrus
PPMI: rostral middle frontal gyrus and superior frontal gyrus
- The observed salient ROIs can be potential **biomarkers** to identify brain disorders.

II. Important Connections

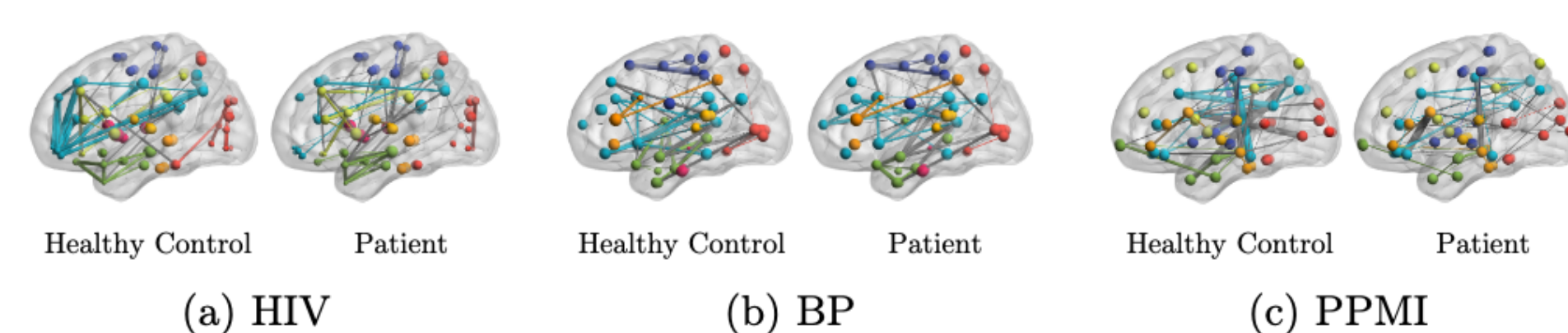


Figure: Important connections on the explanation enhanced brain connection network. Edges connecting nodes within the same neural system (VN, AN, BLN, DMN, SMN, SN, MN, CCN) are colored accordingly.

- HIV: patients excludes rich interactions within the DMN and VN systems
- BP: connections within BLN system of patients are much sparser
- PPMI: connectivity in patients decreases in the SMN and DMN systems

RESOURCES



Paper



Code