

Interpretable Graph Neural Networks for Connectome-Based Brain Disorder Analysis

Hejie Cui¹, Wei Dai¹, Yanqiao Zhu², Xiaoxiao Li³,
Lifang He⁴, and Carl Yang¹

Presenter: Hejie Cui



Paper



Code



EMORY
UNIVERSITY



THE
UNIVERSITY OF
BRITISH
COLUMBIA



LEHIGH
UNIVERSITY

Outline

1. Introduction & Motivation



2. The Proposed Model: IBGNN+

- *The backbone model IBGNN*
- *The globally shared explanation generator*
- *Enhancing the backbone with the learned explanations*

3. Experiments

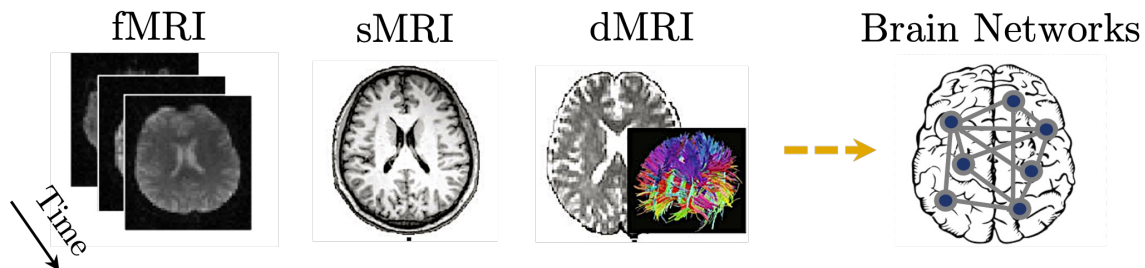
4. Interpretation Analysis

5. Conclusion & Future Work

Introduction: Brain Networks & GNNs

Brain Networks

- Human brains lie at the core of neurobiological systems
- Mapping the connections of the human 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 → universal expressiveness to capture the sophisticated connectome structures
- However, as a family of deep models



Prone to overfitting and lack of **transparency**
and interpretation in predictions!

Introduction: GNN Explanation

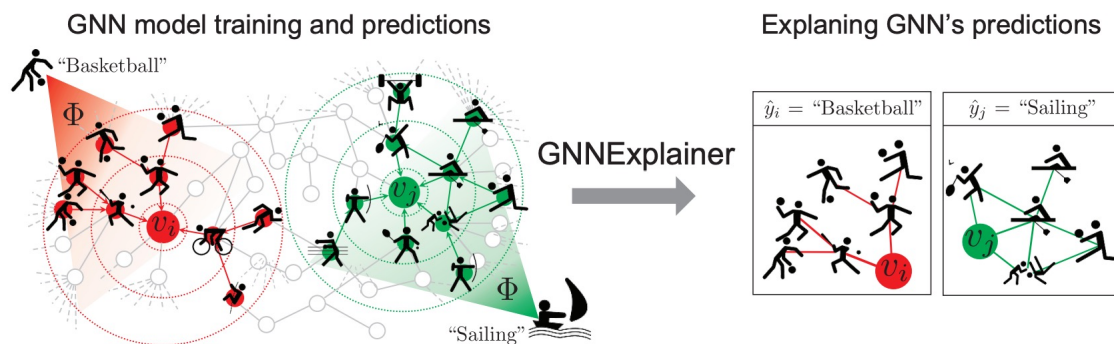
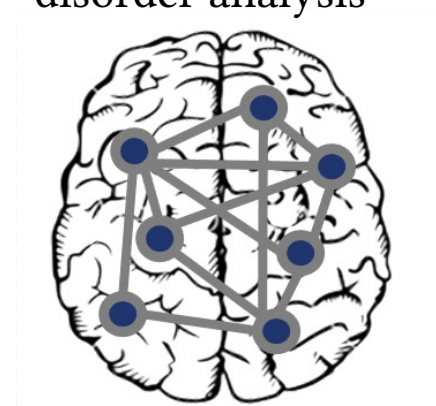
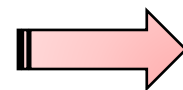


Image source: GNNExplainer: Generating Explanations for Graph Neural Networks

Graph-level connectome-based
disorder analysis



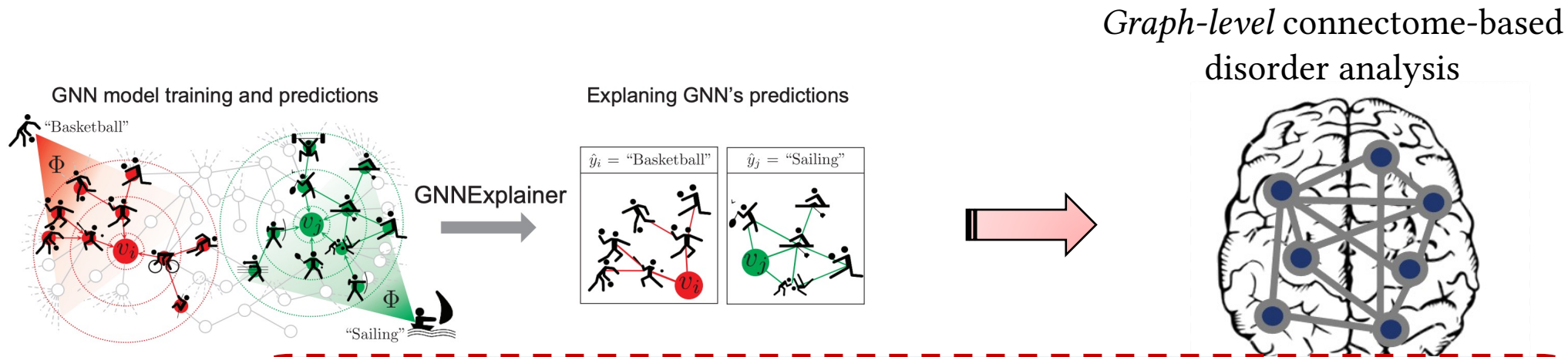
- ❑ Mostly focus on *general graphs* and *node-level* prediction task

- ❑ Produce a *unique explanation* for each subject when applied to graph-level tasks

- ❑ Subjects having the same disorder *share* similar brain network patterns

- ❑ *Unique properties* of brain networks

Introduction: GNN Explanation



❑ Mostly predict

❑ Produce when a

Our motivation:

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

Outline

1. Introduction & Motivation

2. The Proposed Model: IBGNN+

- *The backbone model IBGNN*
- *The globally shared explanation generator*
- *Enhancing the backbone with the learned explanations*

3. Experiments

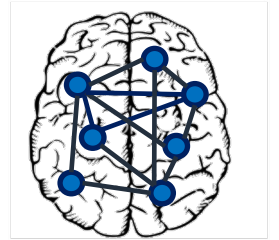
4. Interpretation Analysis

5. Conclusion & Future Work



Problem Definition

- Input: a set of N weighted brain networks, for each network $G = (V, E, W)$
 - $V = \{v_i\}_{i=1}^M$: Regions Of Interest (ROIs) node set of size M
 - $E = V \times V$: edge set of brain connectome
 - $W \in \mathbb{R}^{M \times M}$: weighted adjacency matrix describing the connection strengths between ROIs
- Output:
 - A brain disorder 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



Overview of the Proposed Model: IBGNN+

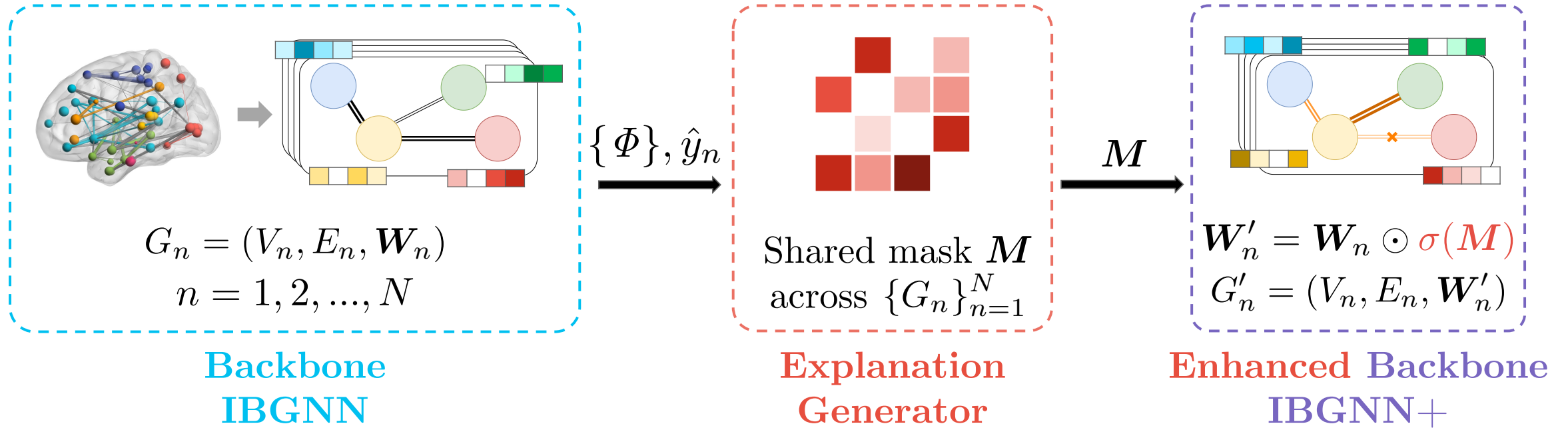


Fig.1: An overview of our proposed framework

- 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 whole model

Module 1: The Backbone Model IBGNN

➤ Edge weights in brain networks can be **both positive and negative values**

Message Vector $\mathbf{m}_{ij} \in \mathbb{R}^D$: concatenate embeddings of a node v_i and its neighbor v_j , and the edge weight w_{ij}

$$\mathbf{m}_{ij}^{(l)} = \text{MLP}_1 \left(\left[\mathbf{h}_i^{(l)}; \mathbf{h}_j^{(l)}; w_{ij} \right] \right)$$

Propagation Rule:

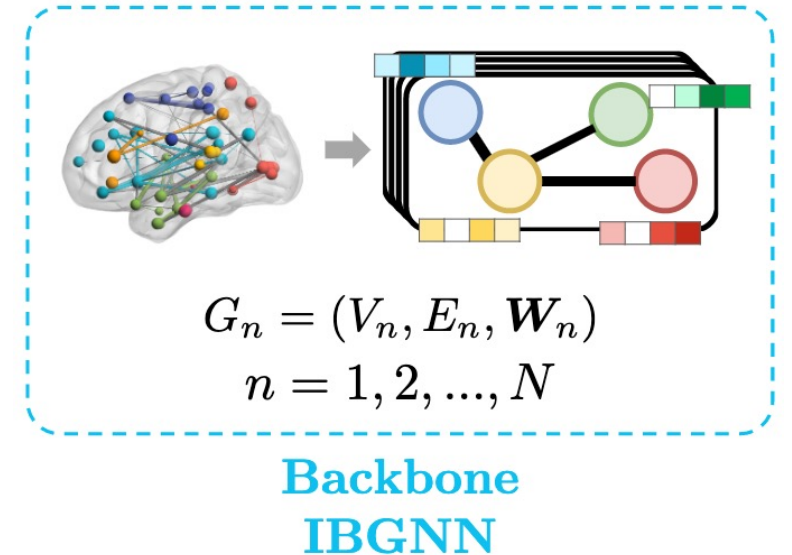
$$\mathbf{h}_i^{(l)} = \xi \left(\sum_{v_j \in \mathcal{N}_i \cup \{v_i\}} \mathbf{m}_{ij}^{(l-1)} \right)$$

Readout Function: summarize all node embeddings to a graph-level one, with MLP and residual connections

$$\mathbf{z} = \sum_{i \in V} \mathbf{h}_i^{(L)} \quad \mathbf{g} = \text{MLP}_2(\mathbf{z}) + \mathbf{z}$$

Training Objectives: supervised 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

- Learn a **globally shared edge mask** $M \in \mathbb{R}^{M \times M}$ that is applied to all brain network subjects in a dataset

Maximize the agreement between the predictions \hat{y} on the original graph G and \hat{y}' on an explanation graph $G' = (V, E, \mathbf{W}')$ induced by a masking matrix \underline{M} , where $\mathbf{W}' = \mathbf{W} \odot \sigma(M)$

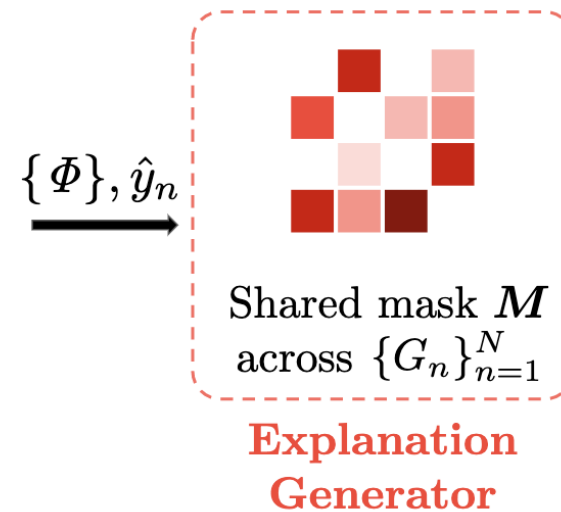
$$\mathcal{L}_{\text{MASK}} = -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C \mathbb{1}[\hat{y}_n = c] \log P_{\Phi}(\hat{y}'_n = \hat{y}_n \mid G'_n)$$

Two regularization terms: encourage the compactness of the explanation and the discreteness of the mask values

$$\mathcal{L}_{\text{SPS}} = \sum_{i,j} M_{i,j} \quad \mathcal{L}_{\text{ENT}} = -(\mathbf{M} \log(\mathbf{M}) + (1 - \mathbf{M}) \log(1 - \mathbf{M}))$$

Training Objectives: weighted sum of supervised cross-entropy \mathcal{L}_{CLF} , agreement loss $\mathcal{L}_{\text{MASK}}$, sparsity loss \mathcal{L}_{SPS} , and discreteness loss \mathcal{L}_{ENT}

$$\mathcal{L} = \mathcal{L}_{\text{CLF}} + \alpha \mathcal{L}_{\text{MASK}} + \beta \mathcal{L}_{\text{SPS}} + \gamma \mathcal{L}_{\text{ENT}}$$



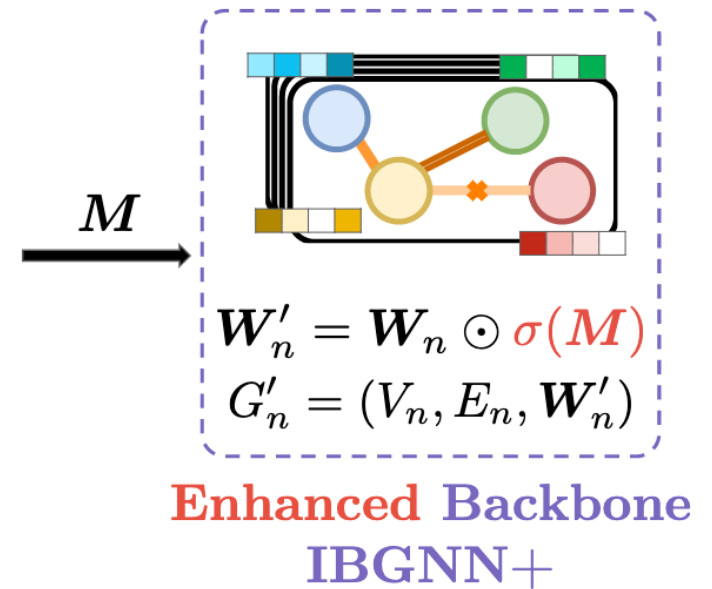
Enhancing the **Backbone** with The **Learned Explanation**: IBGNN+

- Combine the two modules by applying the shared global explanation mask to individual brain networks →

Predictions and Interpretations are produced in a closed-loop for brain network analysis

Suppress random noises

Highlight essential disorder-specific signals



Outline

1. Introduction & Motivation

2. The Proposed Model: IBGNN+

- *The backbone model IBGNN*
- *The globally shared explanation generator*
- *Enhancing the backbone with the learned explanations*

3. Experiments



4. Interpretation Analysis

5. Conclusion & Future Work

Experiments – Datasets

- *Human Immunodeficiency Virus Infection* (HIV)
preprocessed using DPARSF¹ toolbox
- *Bipolar Disorder* (BP)
preprocessed using FSL² toolbox
- *Parkinson's Progression Markers Initiative* (PPMI)
preprocessed using FSL² toolbox and Advanced Normalization Tools (ANTs)³

¹ <http://rfmri.org/DPARSF/>
² <https://fsl.fmrib.ox.ac.uk/fsl/fslwiki/>
³ <http://stnava.github.io/ANTs/>

Dataset	Modality	# samples	Atlas	Size	Response	# classes
HIV	fMRI	70	AAL	90 × 90	Human Immunodeficiency Virus Infection	2
BP	DTI	97	Brodmann	82 × 82	Bipolar Disorder	2
PPMI	DTI	754	Desikan-Killiany	84 × 84	Parkinson's Disorder	2

Tab.1: Dataset summarization

Experiments – Baselines

- Shallow Baselines:
 - M2E¹, MIC², MPCA³, MK-SVM⁴
- Deep Baselines:
 - GCN⁵, GAT⁶, PNA⁷
- SOTA deep models specifically designed for brain networks:
 - BrainNetCNN⁸, BrainGNN⁹

¹ Multi-view multi-graph embedding for brain network clustering analysis. AAAI (2018)

² Clustering on multi-source incomplete data via tensor modeling and factorization. PAKDD (2015)

³ MPCA: Multilinear principal component analysis of tensor objects. IEEE Trans. Neural Netw. (2008)

⁴ Multimodal analysis of functional and structural disconnection in Alzheimer's disease using multiple kernel SVM. Hum. Brain Mapp. (2015)

⁵ Semi-supervised classification with graph convolutional network. ICLR (2017)

⁶ Graph attention networks. ICLR (2018)

⁷ Principal neighbourhood aggregation for graph nets. NeurIPS (2020)

⁸ BrainNetCNN: Convolutional neural networks for brain networks; towards predicting neurodevelopment. NeuroImage (2017)

⁹ BrainGNN: Interpretable brain graph neural network for fMRI analysis. Med. Image Anal. (2021)

Experiments – Prediction Performance

Method	HIV			BP			PPMI		
	Accuracy	F1	AUC	Accuracy	F1	AUC	Accuracy	F1	AUC
M2E	57.14 \pm 19.17	53.71 \pm 19.80	57.50 \pm 18.71	52.56 \pm 13.86	51.65 \pm 13.38	52.42 \pm 13.83	78.69 \pm 1.78	45.81 \pm 4.17	50.39 \pm 2.59
MIC	54.29 \pm 18.95	53.63 \pm 19.44	55.42 \pm 19.10	62.67 \pm 20.92	63.00 \pm 21.61	61.79 \pm 21.74	79.11 \pm 2.16	49.65 \pm 5.10	52.39 \pm 2.94
MPCA	67.14 \pm 20.25	64.28 \pm 23.47	69.17 \pm 20.17	52.56 \pm 13.12	50.43 \pm 14.99	52.42 \pm 13.69	79.15 \pm 0.57	44.18 \pm 0.18	50.00 \pm 0.00
MK-SVM	65.71 \pm 7.00	62.08 \pm 7.49	65.83 \pm 7.41	57.00 \pm 8.89	41.08 \pm 13.44	53.75 \pm 8.00	79.15 \pm 0.57	44.18 \pm 0.18	50.00 \pm 0.00
GCN	70.00 \pm 12.51	68.35 \pm 13.28	73.58 \pm 9.49	55.56 \pm 13.86	50.71 \pm 11.75	61.55 \pm 28.77	78.55 \pm 1.58	47.87 \pm 4.40	59.43 \pm 8.64
GAT	71.43 \pm 11.66	69.79 \pm 10.83	77.17 \pm 9.42	63.34 \pm 9.15	60.42 \pm 7.56	67.07 \pm 5.98	79.02 \pm 1.25	45.85 \pm 3.16	64.40 \pm 6.87
PNA	57.14 \pm 12.78	45.09 \pm 19.62	57.14 \pm 12.78	63.71 \pm 11.34	55.54 \pm 14.06	60.30 \pm 11.89	79.36 \pm 1.84	51.76 \pm 10.32	54.71 \pm 6.77
BrainNetCNN	69.24 \pm 19.04	67.08 \pm 11.11	72.09 \pm 19.01	65.83 \pm 20.64	64.74 \pm 17.42	64.32 \pm 13.72	55.20 \pm 12.63	55.45 \pm 9.15	52.54 \pm 10.21
BrainGNN	74.29 \pm 12.10	73.49 \pm 10.75	75.00 \pm 10.56	68.00 \pm 12.45	62.33 \pm 13.01	74.20 \pm 12.93	69.17 \pm 0.00	44.19 \pm 0.00	45.26 \pm 3.65
IBGNN	82.14 \pm 10.81 [*]	82.02 \pm 10.86 [*]	86.86 \pm 11.65 [*]	73.19 \pm 12.20	72.87 \pm 12.09 [*]	83.64 \pm 9.61 [*]	79.82\pm1.47	51.58 \pm 4.66	70.65 \pm 6.55 [*]
IBGNN+	84.29\pm12.94[*]	83.86\pm13.42[*]	88.57\pm10.89[*]	76.33\pm13.00[*]	76.13\pm13.01[*]	84.61\pm9.08[*]	79.55 \pm 1.67	56.58\pm7.43	72.76\pm6.73[*]

Tab.2: Experiment results (%) on three datasets

- Backbone IBGNN outperforms shallow/deep baselines (up to 11% absolute improvement on BP).
- The explanation enhanced IBGNN+ further improve the backbone by ~9.7% relative improvement.
 - *IBGNN+ can effectively **highlight the disorder-specific signals** while achieving the benefit of **restraining random noises** in individual graphs*

Outline

1. Introduction & Motivation

2. The Proposed Model: IBGNN+

- *The backbone model IBGNN*
- *The globally shared explanation generator*
- *Enhancing the backbone with the learned explanations*

3. Experiments

4. Interpretation Analysis



5. Conclusion & Future Work

Interpretation Analysis – Neural System Mapping

- ROIs on brain networks can be partitioned into different *neural systems*
- 8 commonly used neural systems:
 - Visual Network (VN)
 - Auditory Network (AN)
 - Bilateral Limbic Network (BLN)
 - Default Mode Network (DMN)
 - Somato-Motor Network (SMN)
 - Subcortical Network (SN)
 - Memory Network (MN)
 - Cognitive Control Network (CCN)

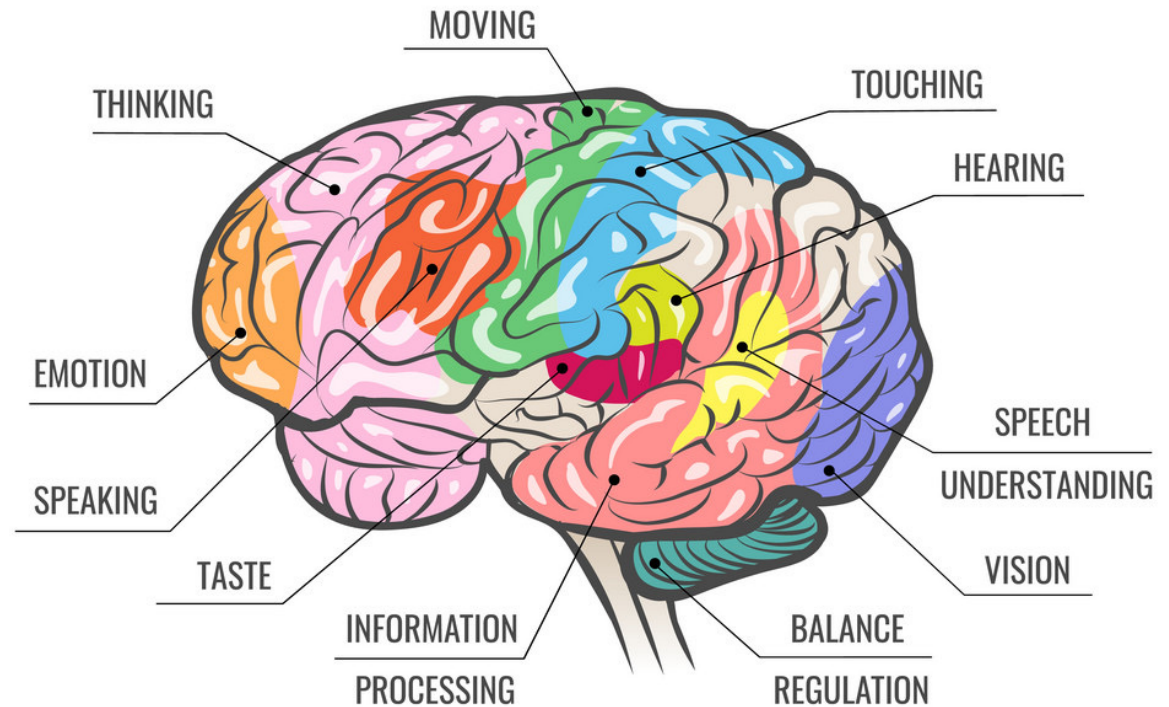


Image source: VectorStock/Human brain function map vector image

Interpretation Analysis – Salient ROIs

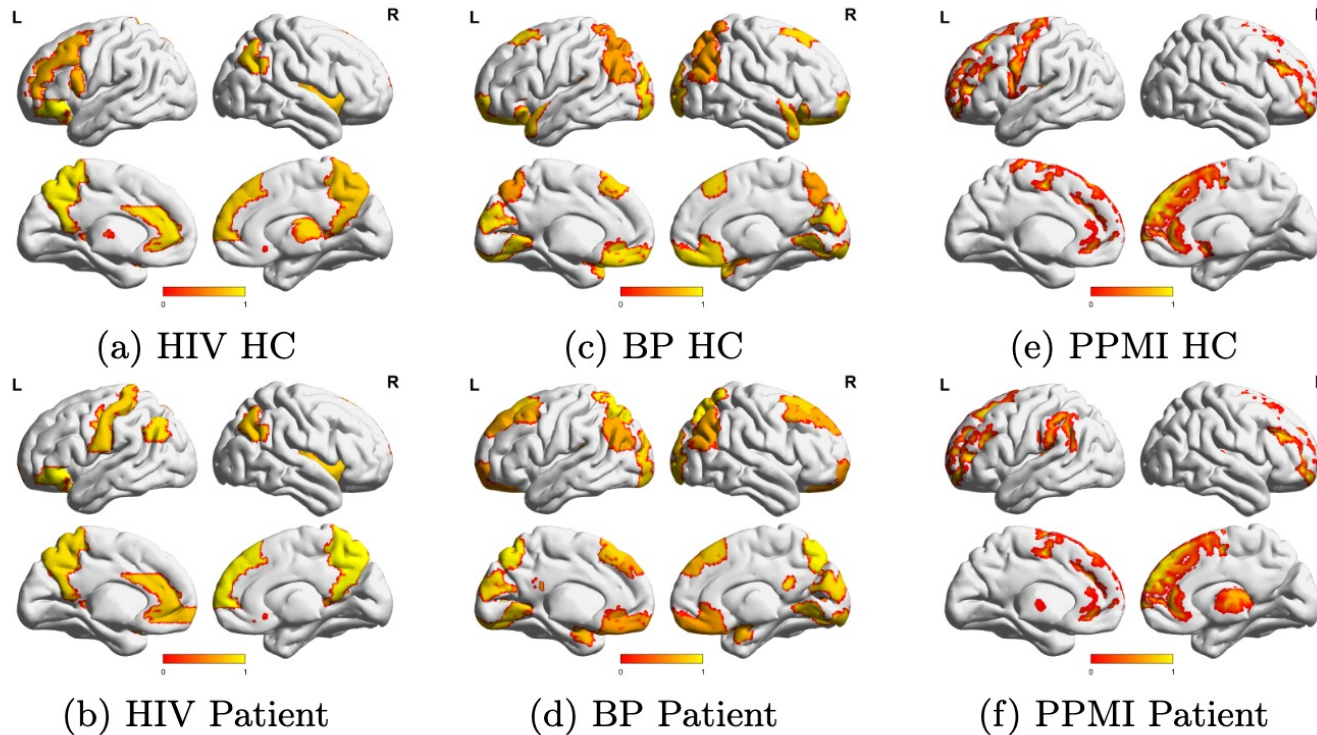


Fig.2: Salient **ROIs** on the explanation enhanced brain connection networks for Health Control (HC) and Patient.

- **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 from each cohort.

Interpretation Analysis – Important Connections

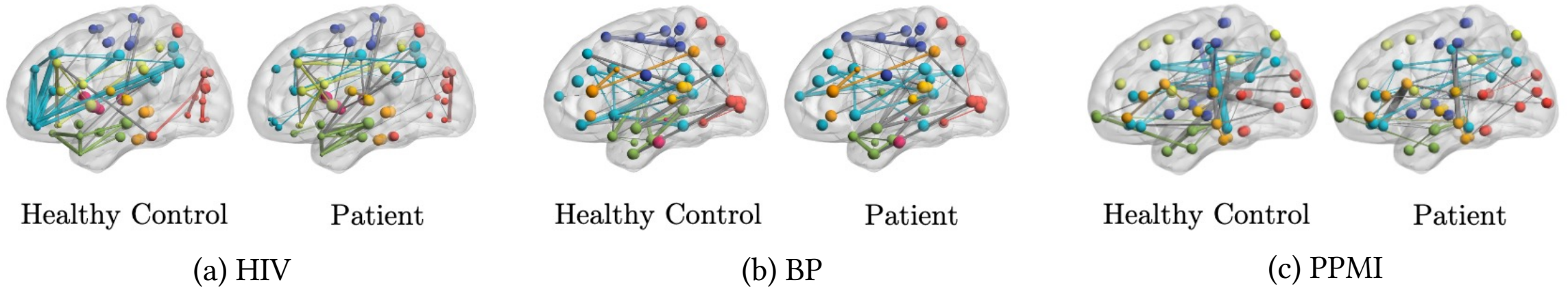


Fig.3: Important **connections** on the explanation enhanced brain connection network. Edges connecting nodes within the same neural system (VN, AN, BLN, DMN, SMN, SN, MN, CCN) are colored accordingly.

Observations:

- **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

Outline

1. Introduction & Motivation

2. The Proposed Model: IBGNN+

- *The backbone model IBGNN*
- *The globally shared explanation generator*
- *Enhancing the backbone with the learned explanations*

3. Experiments

4. Interpretation Analysis

5. Conclusion & Future Work



Summary

Conclusion:

- Propose a novel interpretable GNN framework for connectome-based brain disorder analysis
 - *A brain network-oriented GNN predictor*
 - *A globally shared explanation generator*
- Discover disorder-specific interpretations from the generated explanation mask

Future Work:

- Utilize pre-training and transfer learning techniques to learn across datasets and cohorts

Thanks for Listening!

Presenter: Hejie Cui



Paper



Code