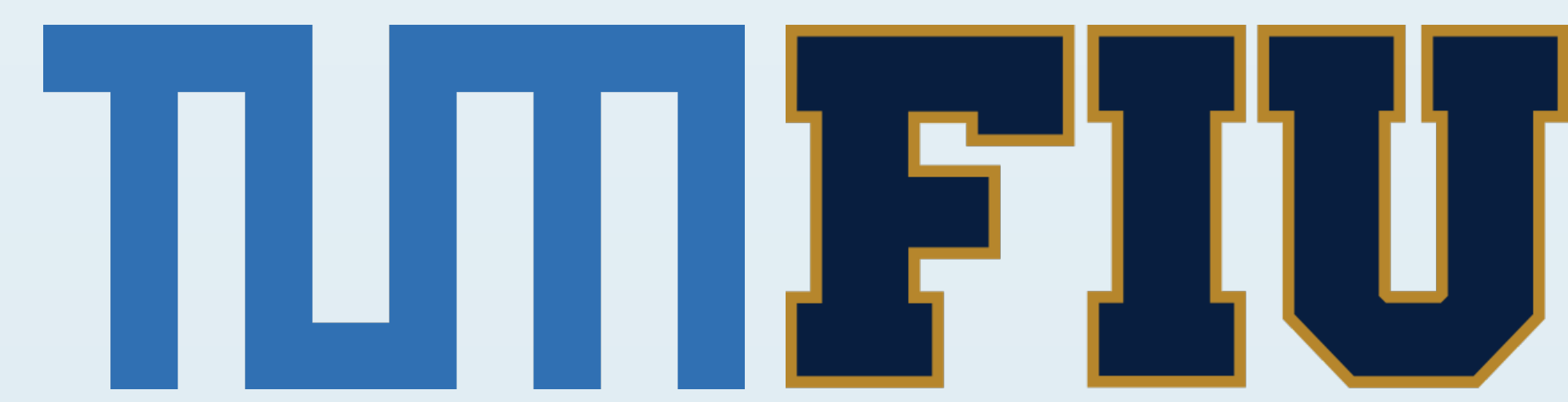


# Factorized Explainer for Graph Neural Networks

Rundong Huang<sup>1</sup>, Farhad Shirani<sup>2\*</sup>, Dongsheng Luo<sup>2\*</sup>

<sup>1</sup> Technical University of Munich, <sup>2</sup> Florida International University



## Objective

- **Post-hoc Instance-Level Explainability** Given model  $f(\cdot)$  and input graph  $G$ , find **minimal** and **sufficient** subgraph  $G_s = \Psi(G)$  w.r.t  $f(\cdot)$ .
- **Research Question 1:** How to quantify **Minimality** and **Sufficiency**?
- **Research Question 2:** How to design explainer mechanisms to optimize the resulting objective?
- **Research Question 3:** How to overcome the locality of the computation graph of GNNs to design explainers for multi-motif scenarios?

## Graph Information Bottleneck

- **Graph Information Bottleneck (GIB):** Tradeoff between **Minimality** and **Sufficiency**

$$\Psi(\cdot) \triangleq \arg \min_{\Psi: G \rightarrow G_s} I(G; G_s) - \alpha I(Y; G_s),$$

- Used to train various explainers such as: PGExplainer, GSAT, and MixupExplainer.
- $I(G; G_s)$  and  $I(Y; G_s)$  difficult to estimate, consequently, surrogates used in practice.

## Shortcomings of GIB

- **Minimality Quantification:**  $I(G; G_s)$  allows leakage of low-entropy components into explanation.
- **Sufficiency Quantification:**  $I(Y; G_s)$  leads to a signaling issue.
- **Theorem 1:** For a graph classification task parametrized by  $P_{G, Y}$ , assume that there exists a mapping  $h: \mathcal{G} \rightarrow \mathcal{Y}$  such that  $G \leftrightarrow h(G) \leftrightarrow Y$  holds. Then, for any  $\alpha > 0$ , there exists an explanation algorithm  $\Psi_\alpha(\cdot)$  such that  $G' \triangleq \Psi_\alpha(G)$  optimizes the objective function in GIB and  $\Psi_\alpha(G) \leftrightarrow h(G) \leftrightarrow G$  holds.

## Modified Objective Function

- **Modified Objective Function:**

$$\Psi(\cdot) = \arg \min_{\Psi: G \rightarrow G_s} \max(\mathbb{E}_G(|G_s|), \beta) + \alpha \mathbb{E}_G(CE(Y; f(G_s))),$$

- $G_s$  is OOD with respect to  $P_G$ .
- Example: In MUTAG, inputs have tens of vertices, however, explanation subgraphs have few vertices.
- The OOD issue is addressed in recent follow-up works.

## Locality of Explanation Method

- **Local Explanation Methods:** Consider a graph classification task  $P_{G, Y}$ , classifier  $f: \mathcal{G} \rightarrow \mathcal{Y}$ , a parameter  $r \in \mathbb{N}$ , and an explanation function  $\Psi: \mathcal{G} \rightarrow \mathcal{G}$ , where  $\mathcal{G}$  is the set of all possible input graphs, and  $\mathcal{Y}$  is the set of output labels. Let  $G' = \Psi(G) = (\mathcal{V}', \mathcal{E}'; \mathbf{Z}', \mathbf{A}')$ . The explanation function  $\Psi(\cdot)$  is called an  $r$ -local explanation function if:

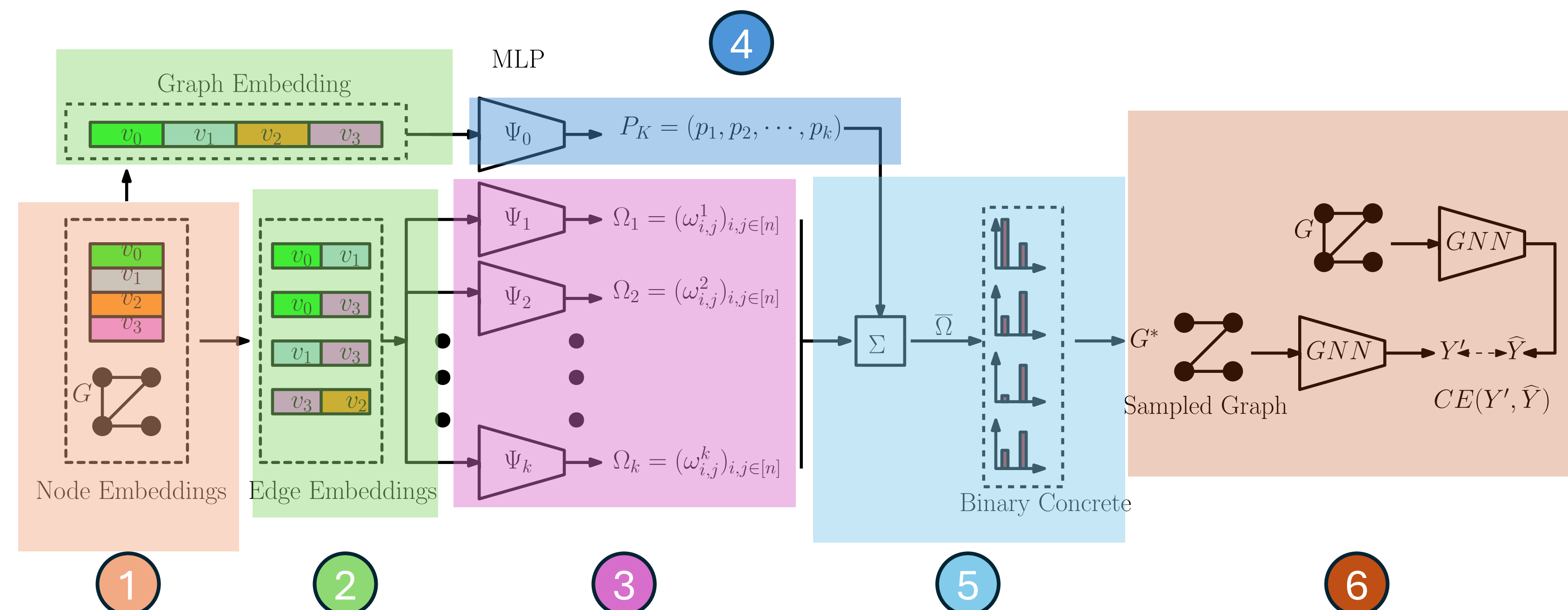
1. The Markov chain  $\mathbb{1}(v \in \mathcal{V}') \leftrightarrow G_{v,r} \leftrightarrow G$  holds for all  $v \in \mathcal{V}$ , where  $\mathbb{1}(\cdot)$  is the indicator function.
2. The edge  $(v, v')$  is in  $\mathcal{E}'$  if and only if  $v, v' \in \mathcal{V}'$  and  $e \in \mathcal{E}$ .

- Explanation methods which rely on GNN node embeddings, such as PGExplainer are local.

## Suboptimality of Local Methods

- **Theorem 2:** Let  $r, r' \in \mathbb{N}$ . There exist classification problems for which:
  - a) The optimal Bayes classification rule  $f^*(g)$  is  $\mathbb{1}(\exists i \in [s]: g_i \subseteq g)$ .
  - b) For any  $r$ -local explanation function, there exists  $\alpha' > 0$  such that the explanation is suboptimal for  $f^*$  in the modified GIB sense for all  $\alpha > \alpha'$  and  $\beta$  equal to maximum number of edges of  $g_i, i \in [s]$ .
  - c) There exists  $k \leq s$ , a parameter  $\alpha' > 0$ , a collection of  $r'$ -local explanation functions  $\Psi_i(\cdot), i \in [k]$ , and an explanation function  $\Psi^*$ , such that for all inputs  $g$ , we have  $\Psi(g) \in \{\Psi_1(g), \Psi_2(g), \dots, \Psi_k(g)\}$  and  $\Psi^*$  is optimal in the modified GIB sense for all  $\alpha > \alpha'$  and  $\gamma$  equal to maximum number of edges of  $g_i, i \in [s]$ .
- The theorem suggest that we can 'patch' together local explainers to construct optimal explanation functions.

## K-FactExplainer



**Step 1:** The original GNN produces node embeddings. **Step 2:** Graph/edge embeddings produced by concatenation. **Step 3:** MLP sequence produce edge probabilities. **Step 4:**  $\Psi_0$  MLP produces weights for each prediction. **Step 5:** Compute weighted average of predictions. **Step 6:** Loss calculated by comparing GNN output for subgraph

## Experimental Results

	BA-Shapes	BA-Community	Tree-Circles	Tree-Grid	BA-2motifs	MUTAG
GRAD	0.882	0.750	0.905	0.667	0.717	0.783
ATT	0.815	0.739	0.824	0.612	0.674	0.765
RGExp.	0.985 $\pm$ 0.013	0.919 $\pm$ 0.017	0.787 $\pm$ 0.099	<b>0.927</b> $\pm$ 0.032	0.657 $\pm$ 0.107	0.873 $\pm$ 0.028
DEGREE	0.993 $\pm$ 0.005	0.957 $\pm$ 0.010	<b>0.902</b> $\pm$ 0.022	0.925 $\pm$ 0.040	0.755 $\pm$ 0.135	0.773 $\pm$ 0.029
GNNExp.	0.742 $\pm$ 0.006	0.708 $\pm$ 0.004	0.540 $\pm$ 0.017	0.714 $\pm$ 0.002	0.499 $\pm$ 0.004	0.606 $\pm$ 0.003
PGExp.	0.999 $\pm$ 0.000	0.825 $\pm$ 0.040	0.760 $\pm$ 0.014	0.679 $\pm$ 0.008	0.566 $\pm$ 0.004	0.843 $\pm$ 0.162
K-FactExplainer	<b>1.000</b> $\pm$ 0.000	<b>0.974</b> $\pm$ 0.004	0.779 $\pm$ 0.004	0.770 $\pm$ 0.004	<b>0.821</b> $\pm$ 0.005	<b>0.915</b> $\pm$ 0.010

Explanation faithfulness in terms of AUC-ROC on edges under six datasets.

## Acknowledgements

This project was partially supported by NSF grants IIS-2331908 and CCF-2241057. The views and conclusions contained in this paper are those of the authors and should not be interpreted as representing any funding agencies.