



NEURAL INFORMATION
PROCESSING SYSTEMS



Microsoft

Train Once and Explain Everywhere: Pre-training Interpretable Graph Neural Networks

Jun Yin*, Chaozhuo Li*, Hao Yan, Jianxun Lian, Senzhang Wang

Central South University, Changsha, China

Microsoft Research Asia, Beijing, China

NeurIPS 2023

CONTENT



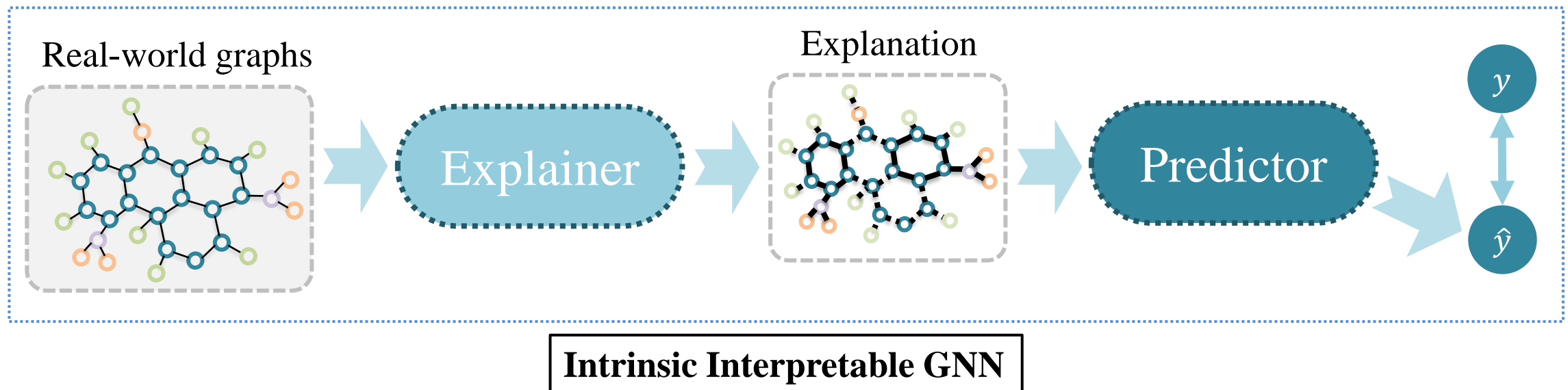
01 BACKGROUND

02 MODEL FRAMEWORK

03 EXPERIMENT

04 SUMMARY

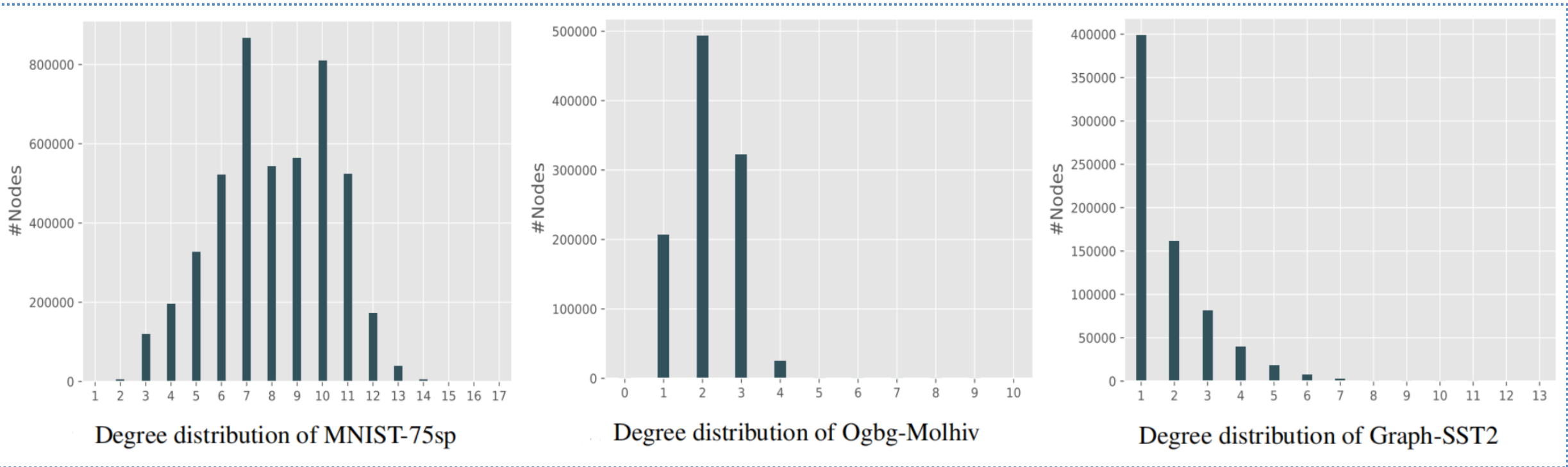
- GNNs have achieved remarkable success in various applications.
- Black-box nature makes it hard to understand the inner decision-making mechanism.
- Intrinsic interpretable GNNs aim to provide transparent predictions by identifying the influential fraction of the input graph that guides the model prediction.



BACKGROUND



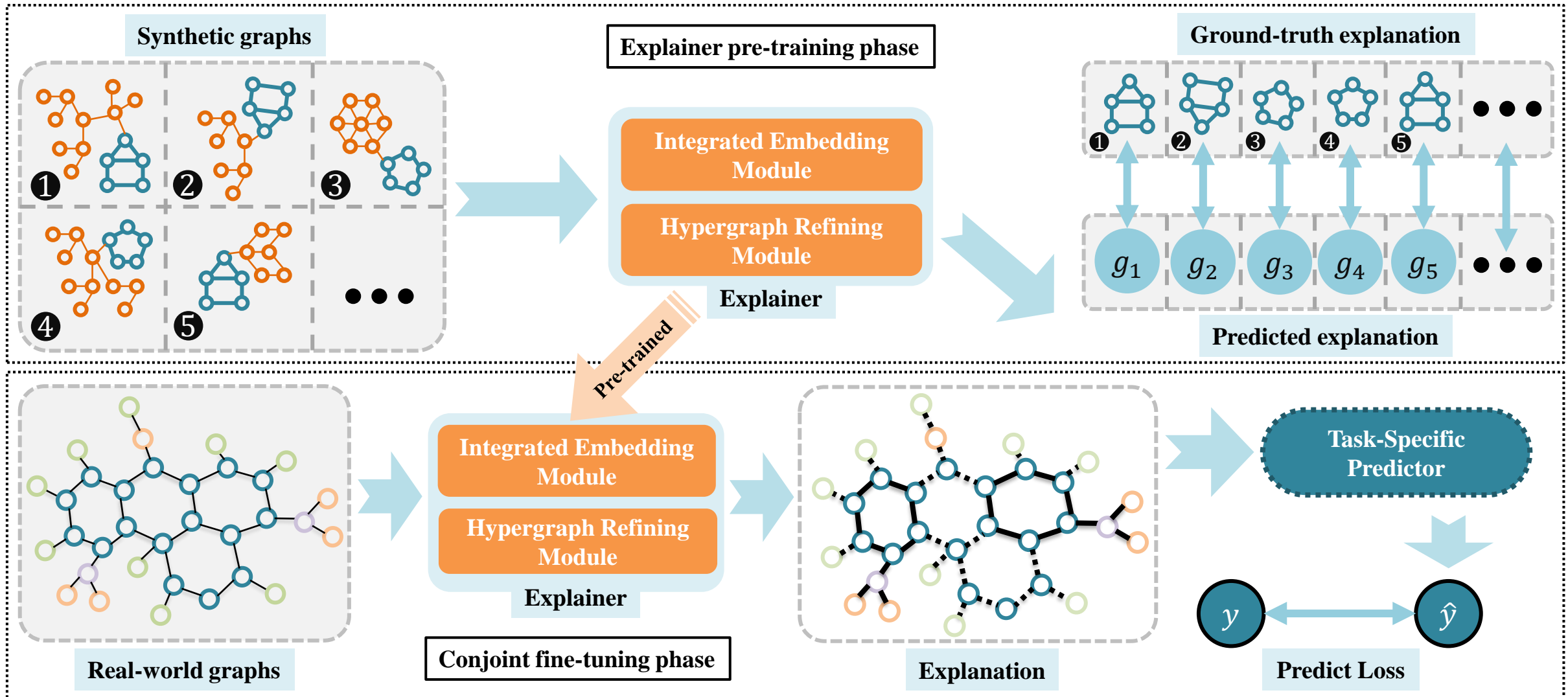
- Existing GNN explanation methods are dataset-specific.
- How to construct a interpretable GNN that can generalize to different datasets?



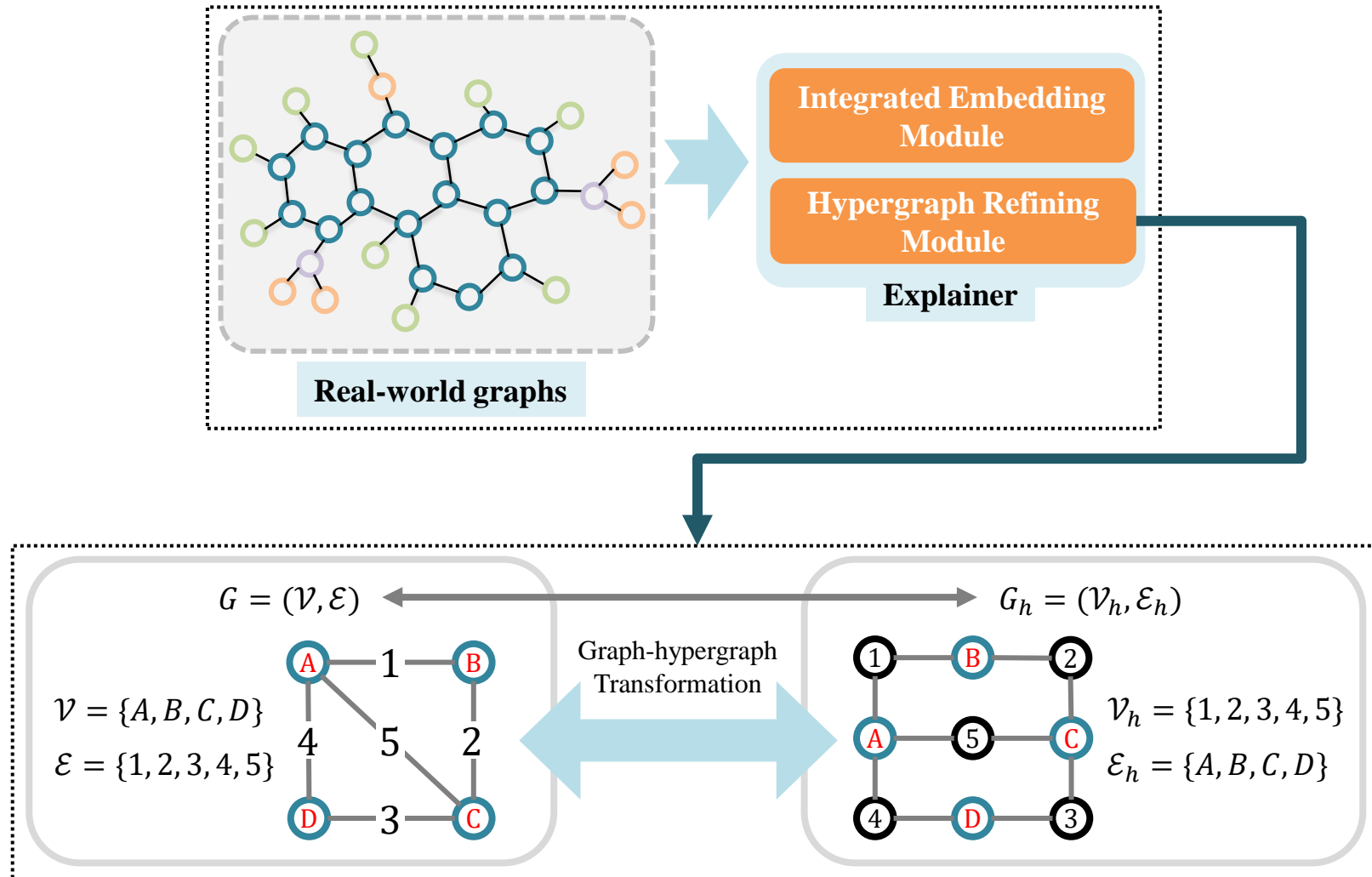
MODEL FRAMEWORK



- Pre-training on synthetic graphs.



- Hyper-graph refining module.



- Interpretation Performance.

Table 1: Interpretation Performance (ROC-AUC) Comparison. The underlined results highlight the best baselines. The **bold** results mean the π -GNN or π -GNN_{DFT} outperform the best baselines.

Model	BA-2Motifs	Mutag	MNIST-75sp	Spurious-Motif		
				$b = 0.5$	$b = 0.7$	$b = 0.9$
GNExplainer	67.35 ± 3.29	61.98 ± 5.45	59.01 ± 2.04	62.62 ± 1.35	62.25 ± 3.61	58.86 ± 1.93
PGExplainer	84.59 ± 9.09	60.91 ± 17.10	69.34 ± 4.32	69.54 ± 5.64	72.33 ± 9.18	72.34 ± 2.91
GraphMask	92.54 ± 8.07	62.23 ± 9.01	73.10 ± 6.41	72.06 ± 5.58	73.06 ± 4.91	66.68 ± 6.96
IB-Subgraph	86.06 ± 28.37	91.04 ± 6.59	51.20 ± 5.12	57.29 ± 14.35	62.89 ± 15.59	47.29 ± 13.39
DIR	82.78 ± 10.97	64.44 ± 28.81	32.35 ± 9.39	78.15 ± 1.32	77.68 ± 1.22	49.08 ± 3.66
GIN-GSAT	<u>98.74</u> ± 0.55	<u>99.60</u> ± 0.51	83.36 ± 1.02	78.45 ± 3.12	74.07 ± 5.28	71.97 ± 4.41
PNA-GSAT	93.77 ± 3.90	99.07 ± 0.50	<u>84.68</u> ± 1.06	<u>83.34</u> ± 2.17	<u>86.94</u> ± 4.05	<u>88.66</u> ± 2.44
π -GNN	99.33 ± 0.63	99.81 ± 0.17	92.77 ± 0.80	93.24 ± 0.72	96.92 ± 0.85	96.39 ± 0.92
π -GNN _{DFT}	93.19 ± 1.48	95.29 ± 0.67	85.18 ± 1.08	86.29 ± 2.22	87.43 ± 2.47	89.64 ± 2.26

- Prediction Performance.

Table 2: Prediction Performance (Acc) Comparison. The underlined results highlight the best baselines. The **bold** results mean the π -GNN or π -GNN_{DFT} outperform the best baselines.

Model	Molhiv(AUC)	Graph-SST2	MNIST-75sp	Spurious-Motif		
				$b = 0.5$	$b = 0.7$	$b = 0.9$
GIN	76.69 ± 1.25	82.73 ± 0.77	95.74 ± 0.36	39.87 ± 1.30	39.04 ± 1.62	38.57 ± 2.31
PNA	78.91 ± 1.04	79.87 ± 1.02	87.20 ± 5.61	68.15 ± 2.39	<u>66.35</u> ± 3.34	<u>61.40</u> ± 3.56
IB-Subgraph	76.43 ± 2.65	<u>82.99</u> ± 0.67	93.10 ± 1.32	54.36 ± 7.09	48.51 ± 5.76	46.19 ± 5.63
DIR	76.34 ± 1.01	82.32 ± 0.85	88.51 ± 2.57	45.49 ± 3.81	41.13 ± 2.62	37.61 ± 2.02
GIN-GSAT	76.47 ± 1.53	82.95 ± 0.58	<u>96.24</u> ± 0.17	52.74 ± 4.08	49.12 ± 3.29	44.22 ± 5.57
PNA-GSAT	<u>80.24</u> ± 0.73	80.92 ± 0.66	<u>93.96</u> ± 0.92	<u>68.74</u> ± 2.24	64.38 ± 3.20	57.01 ± 2.95
π -GNN	80.86 ± 0.61	88.05 ± 0.43	96.89 ± 0.20	74.67 ± 0.63	77.52 ± 0.77	77.46 ± 0.96
π -GNN _{DFT}	79.71 ± 1.08	83.48 ± 1.20	92.89 ± 0.95	70.78 ± 1.63	71.02 ± 1.43	72.61 ± 1.75

- An interpretable GNN that can generalize to different graph datasets.
- Synthetic pre-training process.
- Hyper-graph transformation based edge representation.



Microsoft

A decorative graphic on the left side of the slide, composed of purple, organic, branching shapes that resemble neural connections or a stylized tree.

NEURAL INFORMATION PROCESSING SYSTEMS

THANKS