



NExT++ Research Center

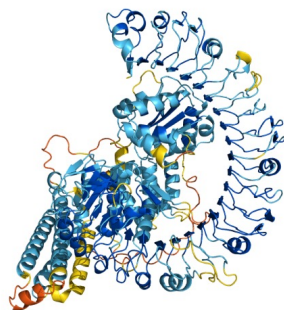
A Joint Research Collaboration Between NUS & Tsinghua University

---

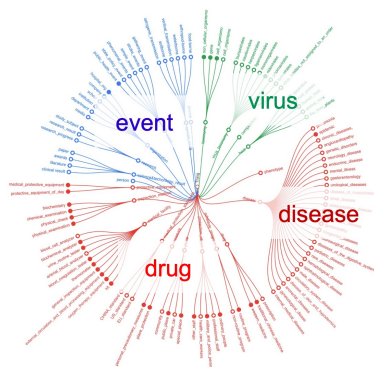


# Towards **Multi-Grained Explainability** for **Graph Neural Networks**

Xiang Wang, Yingxin Wu, An Zhang, Xiangnan He, Tat-Seng Chua



Protein Structure



COVID Graph



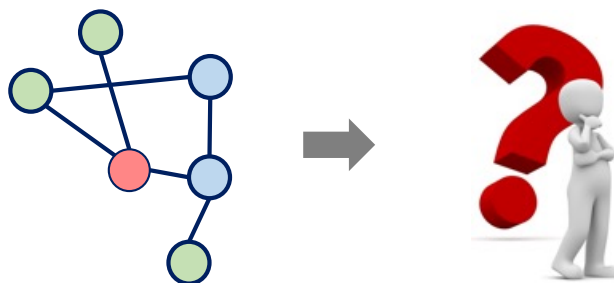
Knowledge Graph



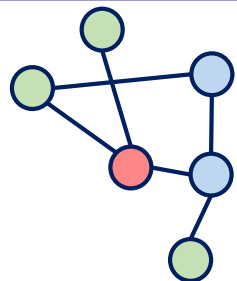
Social Network

## Tasks of Graph Learning:

- Node Classification
- Graph Classification
- Link Prediction ...



## Core of GNNs



Input Graph  $G$



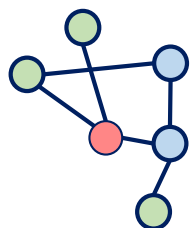
GNN Model  $f$



Output Prediction  $\hat{y}$

- Graph Structure Guides Representation Learning
- Information Propagation & Aggregation

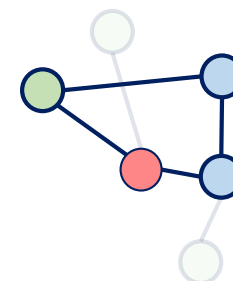
## Explainability of GNNs



Input Graph  $G$



Which fraction of the input graph is **most influential** to the model's decision?

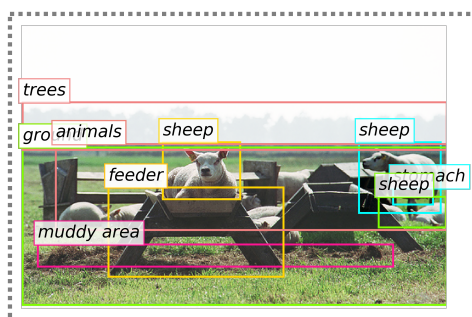


Explanatory Subgraph  $G_s$

Output Prediction  $\hat{y}$

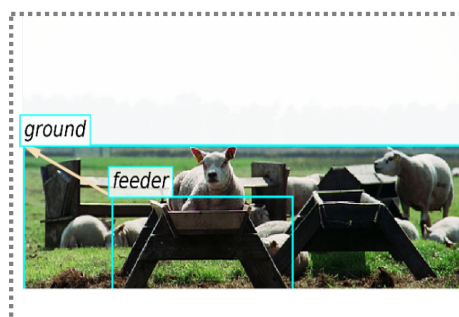
Case: Scene graph classification

Question: Which fraction of the input graph is most influential to the model's decision?



Input Scene Graph

Output Prediction  
 $\hat{y}$ : Farm



Explanation 1



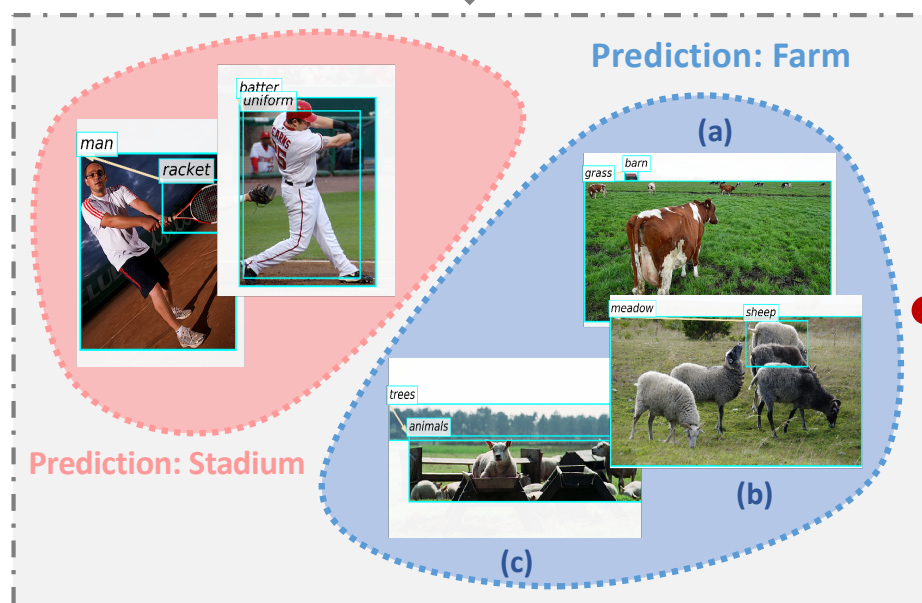
Explanation 2

	Related works	Potential drawback
<b>Local Explainability:</b> Interprets each instance independently.	GNExplainer [Ying et al. 2019] PGM-Explainer [Vu et al. 2020]	They hardly exhibit the class-wise patterns
<b>Global Explainability:</b> Systematizes the globally important patterns.	PGExplainer [Luo et al. 2020] XGNN [Yuan et al. 2020]	They might be trivial in the local context.

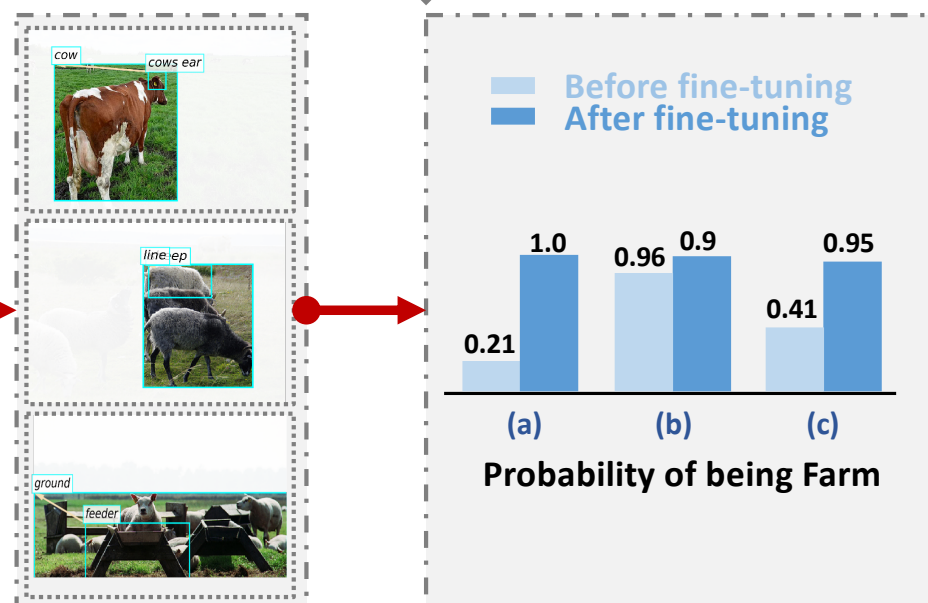
<https://github.com/Wuyxin/ReFine>.

What **class-wise knowledge** does the GNN leverage to make predictions **in general**?

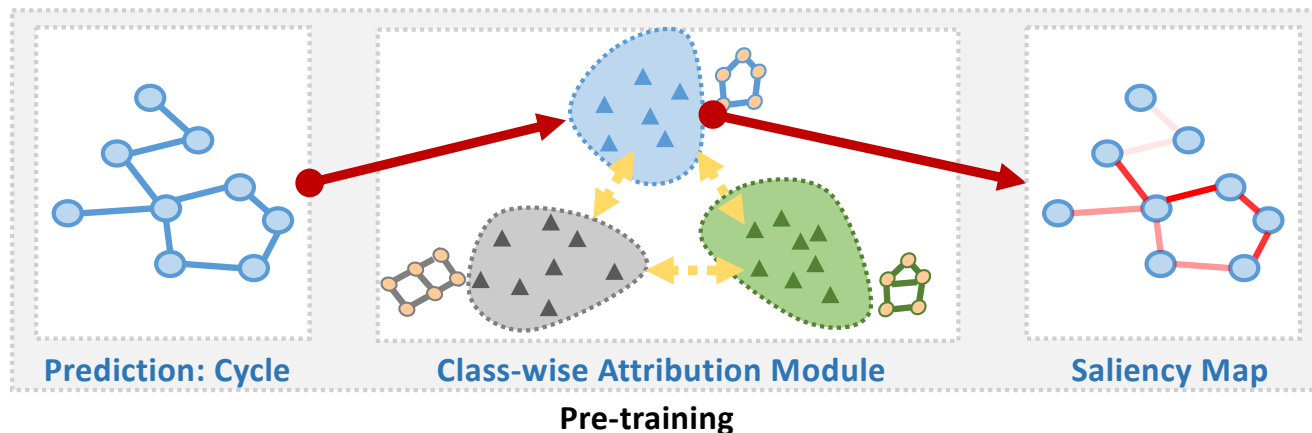
Why the GNN model made the **certain prediction** for the instance **at hand**?



Pre-training towards Global Explainability



Fine-Tuning towards Local Explainability

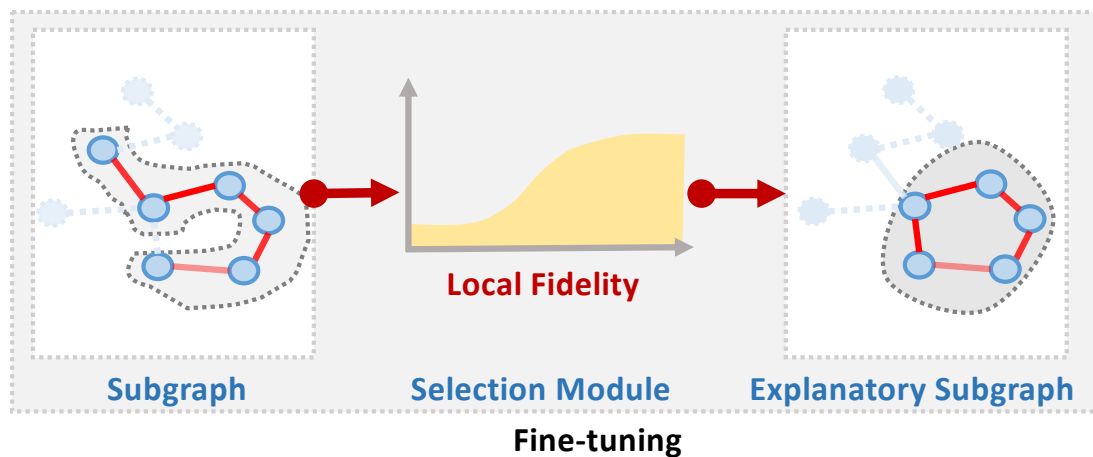


$$\mathcal{L}_1 = MI(Y, \mathbf{M} \odot G_{att})$$

Negative mutual information

$$\min_{\theta} \mathcal{L}_1 + \gamma \mathcal{L}_{cts}$$

Contrastive Loss



$$\theta'_0 = \theta$$

$$\min_{\theta'} \mathcal{L}_2 = MI(Y, G_{exp})$$

where  $G_{exp} = \mathbf{Top}_\rho(G_{att})$

User-defined ratio

Table 1: Comparison of our ReFine and other baseline explainers

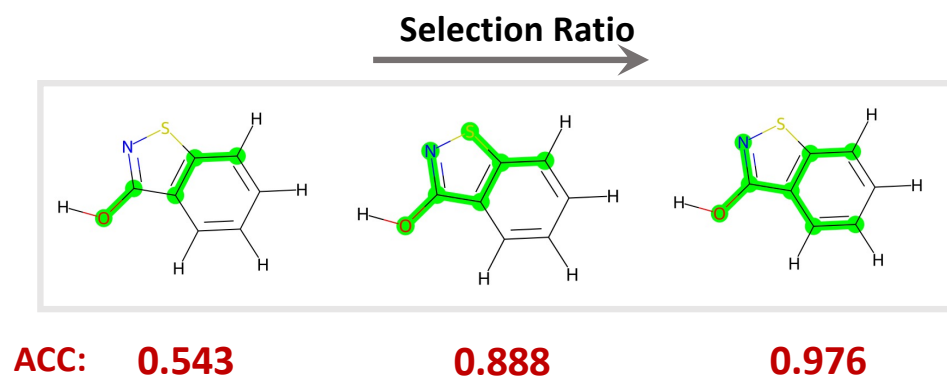
	Mutagenicity	VG-5	MNIST	BA-3motif	
	ACC-AUC	ACC-AUC	ACC-AUC	ACC-AUC	Recall@5
SA[9]	0.769	0.769	0.559	0.518	0.243
GNNExplainer[6]	0.895±0.010	0.895±0.003	0.535±0.013	0.528±0.005	0.157±0.002
PG-Explainer[7]	0.631±0.008	0.790±0.004	0.504±0.010	0.586±0.004	0.293±0.001
PGM-Explainer[19]	0.714±0.007	0.792±0.001	0.615±0.003	0.575±0.002	0.250±0.000
<b>ReFine-CT</b>	0.888±0.008	0.891±0.002	0.526±0.007	0.610±0.004	0.248±0.001
<b>ReFine-FT</b>	0.945±0.011	0.906±0.002	0.587±0.008	0.616±0.003	0.299±0.002
<b>ReFine</b>	<b>0.955±0.005</b>	<b>0.914±0.001</b>	<b>0.636±0.003</b>	<b>0.630±0.006</b>	<b>0.304±0.000</b>
Improvement	6.7%	2.1%	3.4%	7.5%	3.8%

Table 2: Performance under different selection ratios before and after fine-tuning.

ACC@ $\rho$	Mutagenicity		VG-5		MNIST		BA-3motif	
	0.4	0.6	0.4	0.6	0.4	0.6	0.4	0.6
ReFine-FT	96.8%	94.0%	91.3%	91.4%	41.4%	61.4%	36.0%	65.7%
ReFine	97.8%	96.2%	92.2%	93.4%	71.4%	82.0%	39.0%	72.8%
Improvement	+1.0%	+2.2%	+0.9%	+2.0%	+30.0%	+20.6%	+3.0%	+7.1%

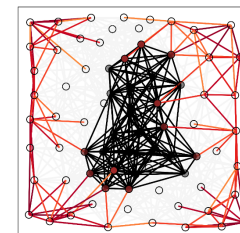
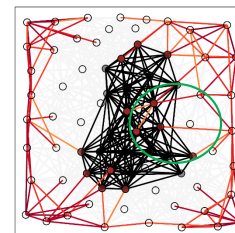
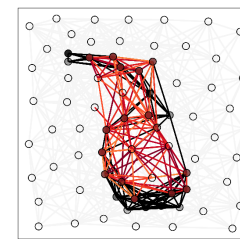
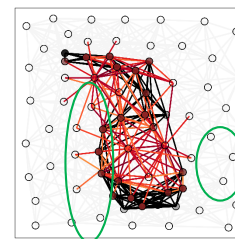
<https://github.com/Wuyxin/ReFine>.





**Pre-trained**

**Fine-tuned**







## Summary



- Local explainability & Global explainability present different views of GNN models.
- Multi-grained explainability can offer more reliable & faithful explanations.

**Check out our code and models at**

- <https://github.com/Wuyxin/ReFine>.

# THANK YOU!