

Bagel: A Benchmark for Assessing Graph Neural Network Explanations

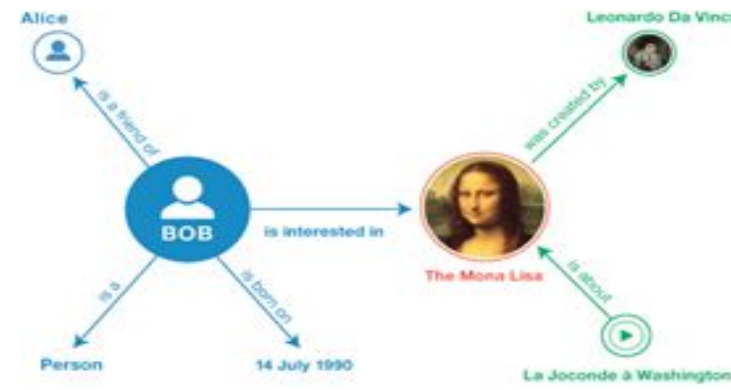
Mandeep Rathee, Thorben Funke, Avishek Anand, and Megha Khosla



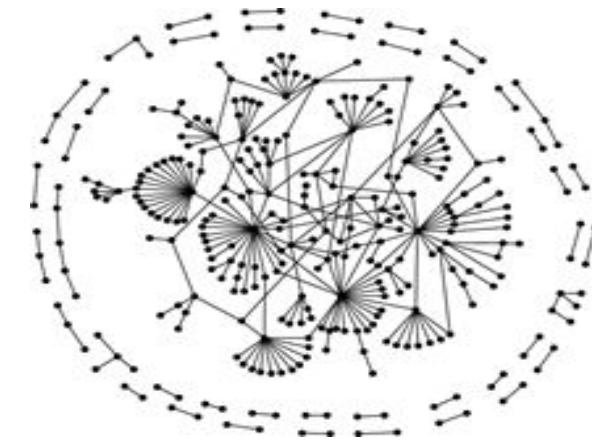
Graphs Are Everywhere



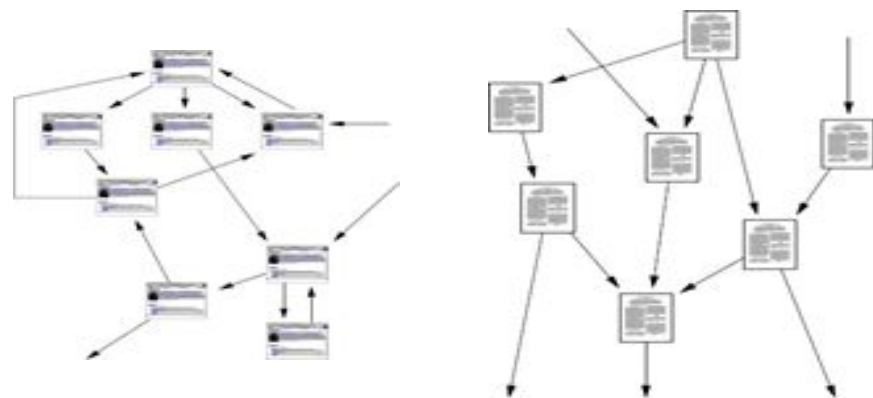
Social Networks



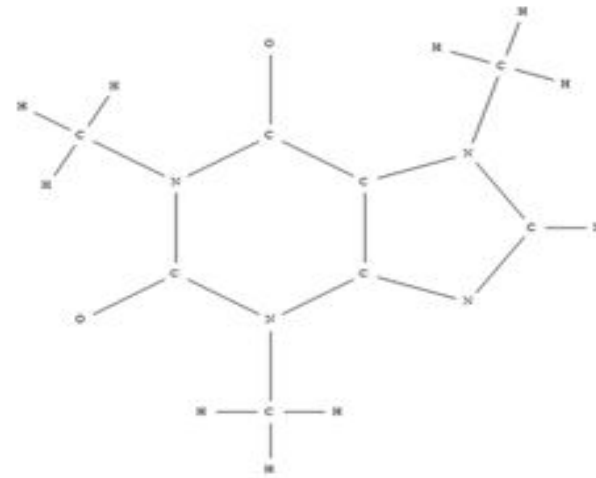
Knowledge Graphs



Biological Networks

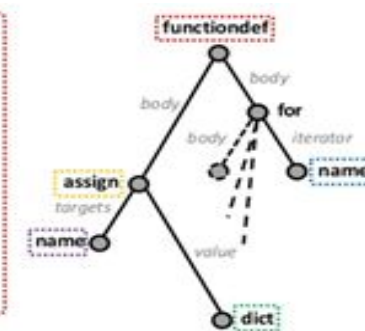


Citation Networks



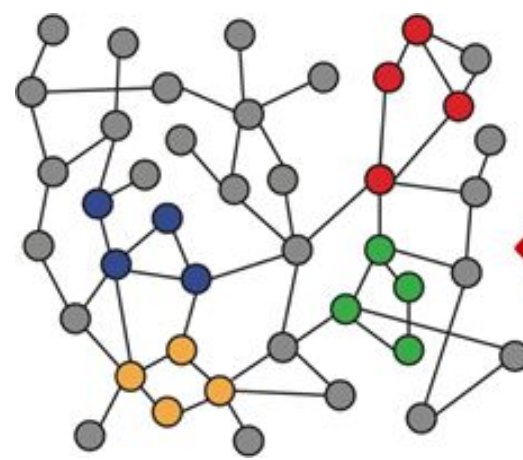
Molecules

```
def encode(obj):  
    """  
    Encode a (possibly nested)  
    dictionary containing complex values  
    into a form that can be serialized  
    using JSON.  
    """  
    e = {}  
    for key, value in obj.items():  
        if isinstance(value, dict):  
            e[key] = encode(value)  
        elif isinstance(value, complex):  
            e[key] = {'type': 'complex',  
                    'r': value.real,  
                    'i': value.imag}  
    return e  
  
import ast  
tree = ast.parse("...")
```

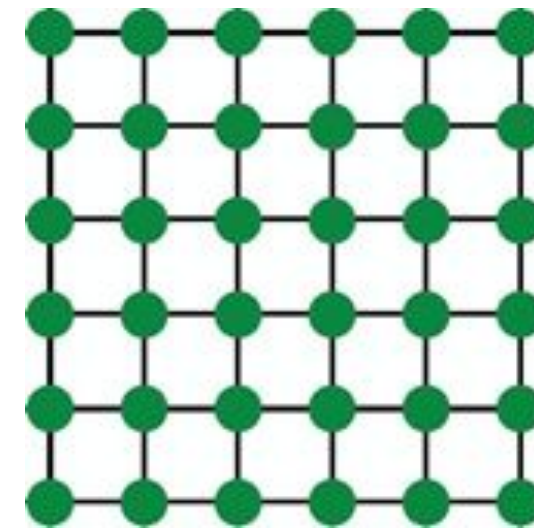
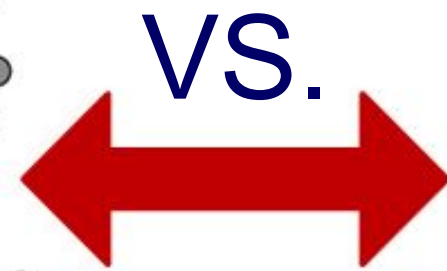


Code

Graphs Are Complex



Random Network



Images

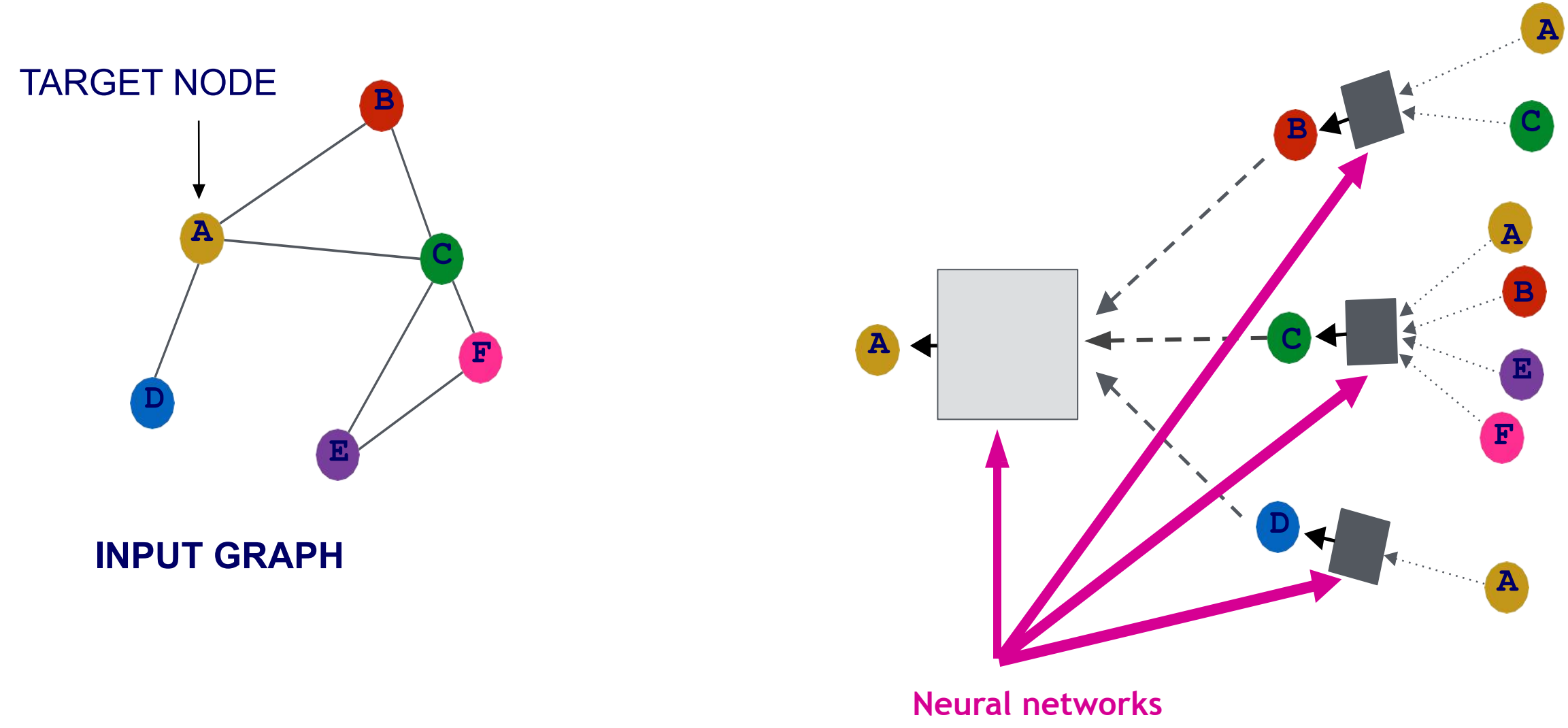


Text/Audio



Tabular

Graph Neural Networks (GNNs)



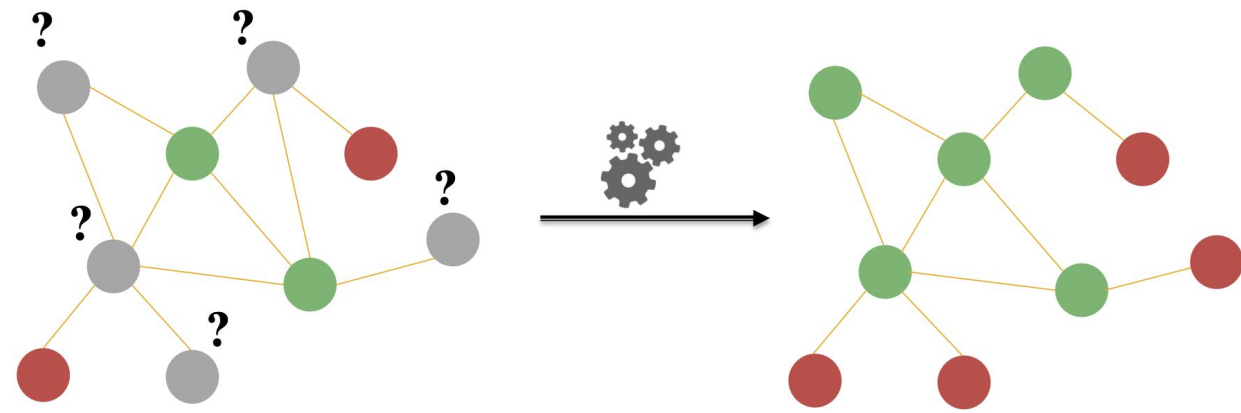
Each node defines a computation graph

- Each edge in this graph is a transformation/aggregation function

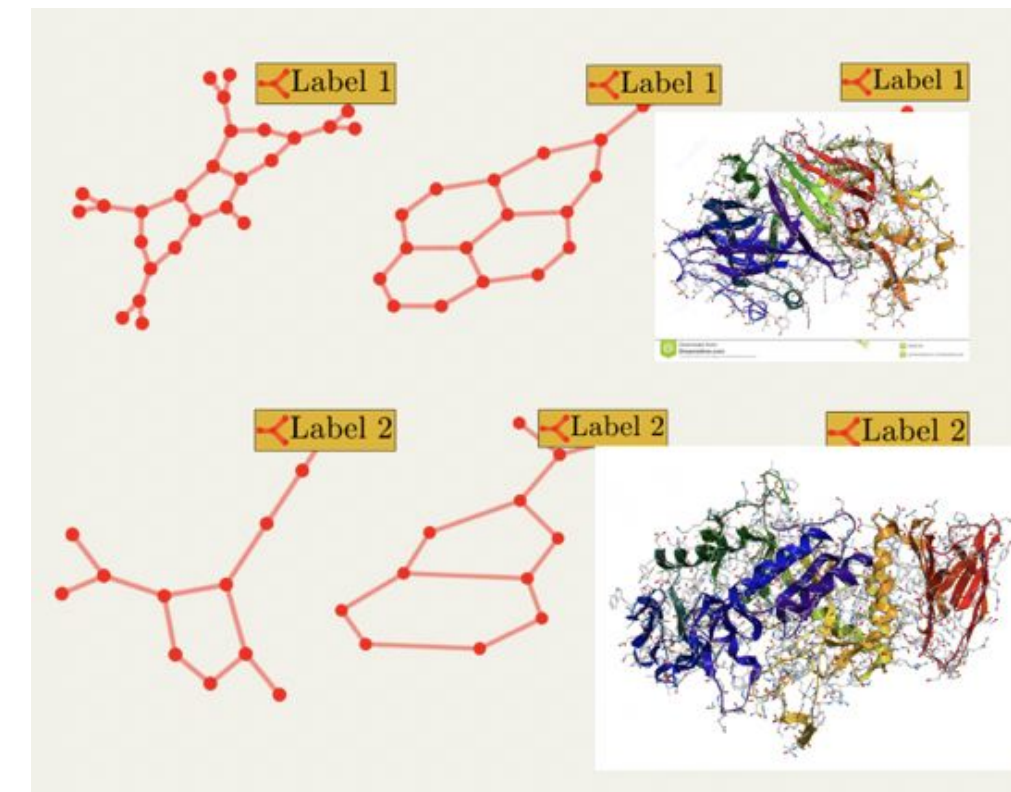
e.g., GCN and GAT

Scarselli et al. 2005. [The Graph Neural Network Model](#). *IEEE Transactions on Neural Networks*.

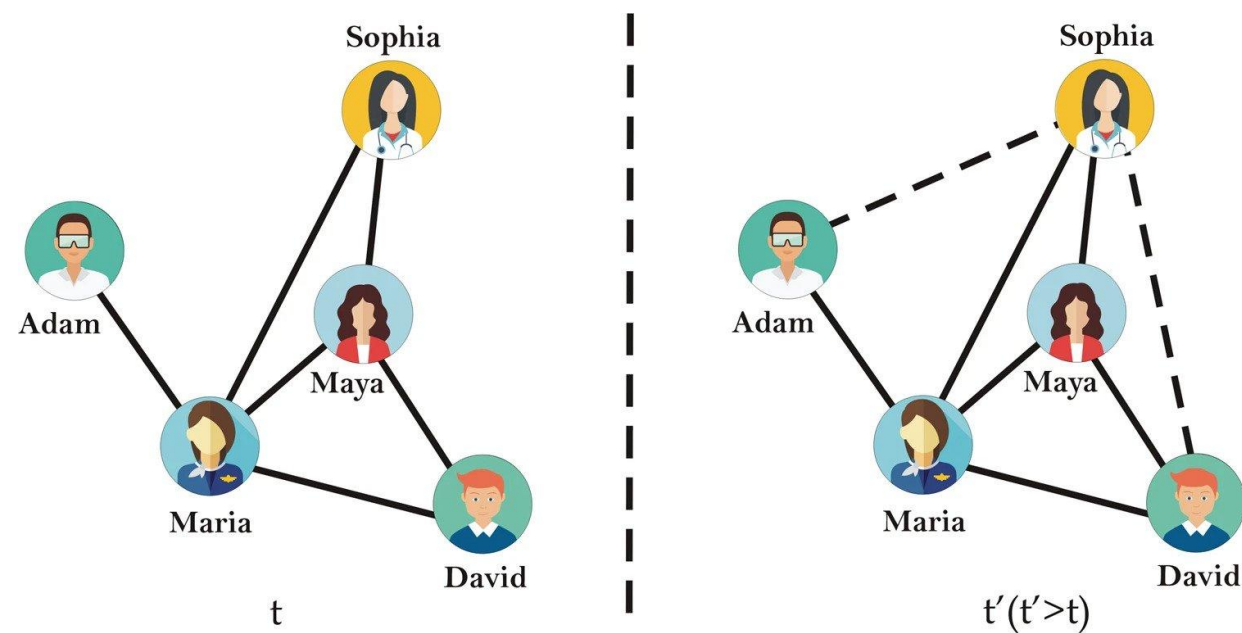
Tasks on Graph-Structured Data



Node Classification



Graph Classification



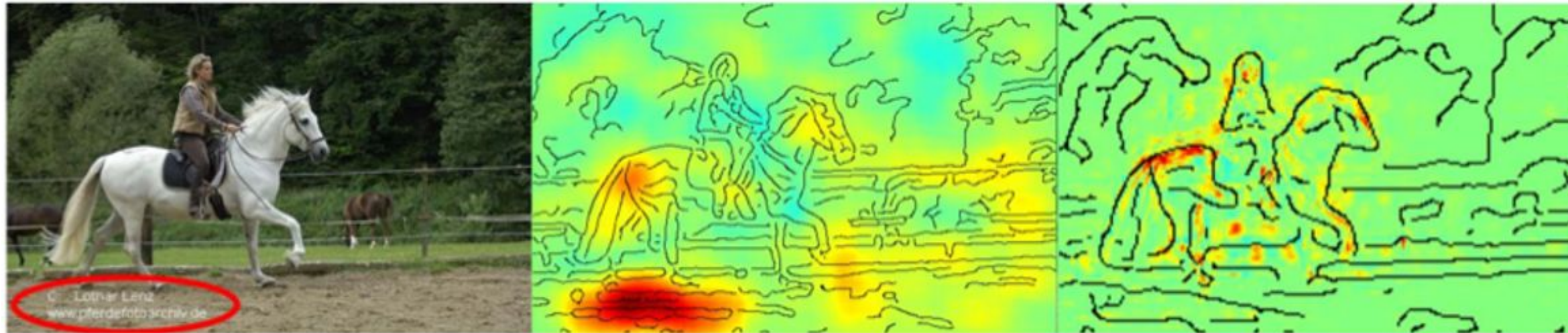
Link Prediction

Applications in

- Health
- Recommendation
- Finance

Why Explainability?

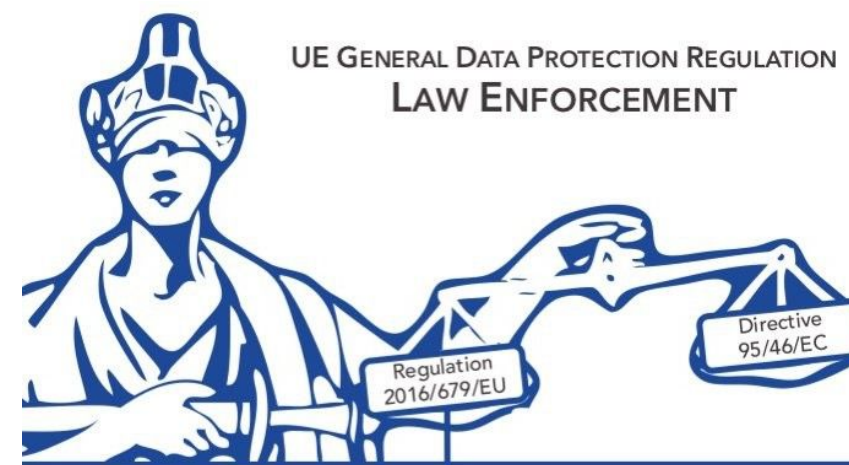
Given a machine learning model and machine learning task, we are interested in finding the rationales behind the prediction.



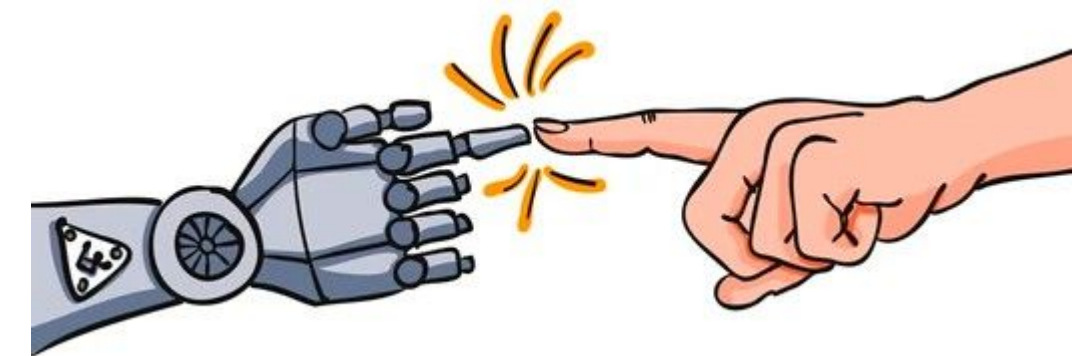
Right for the Right reasons



Utilize insights to improve models

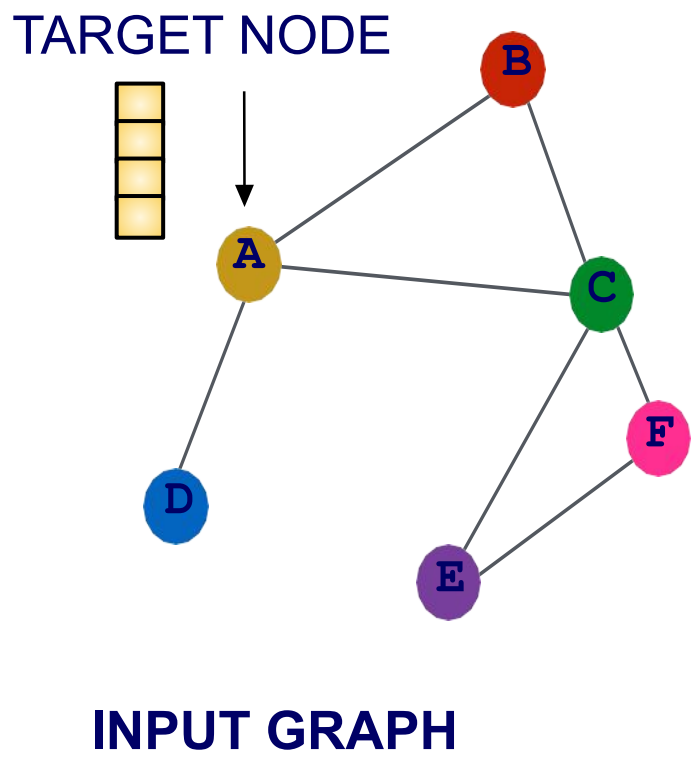


Legal recourse



Improve trust

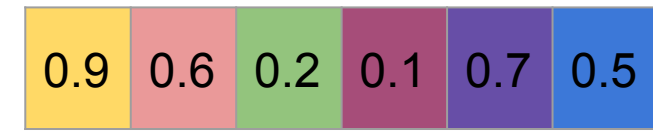
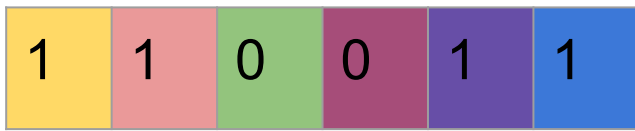
Explainability in GNNs



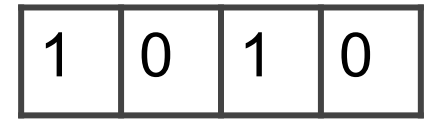
Hard Mask

Soft Mask

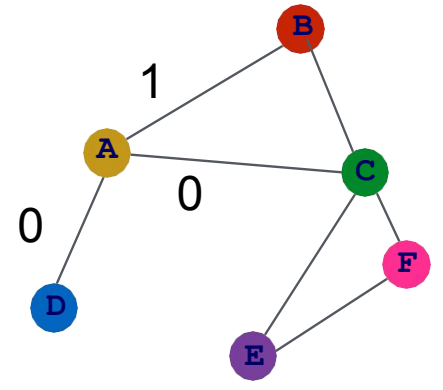
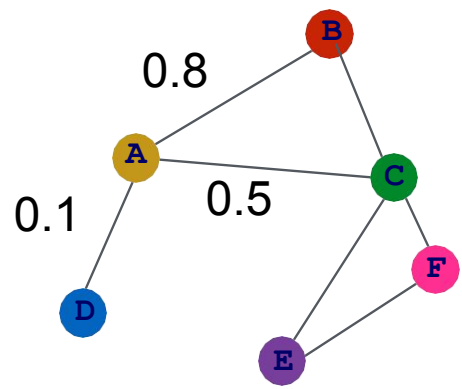
Node Level



Feature Level

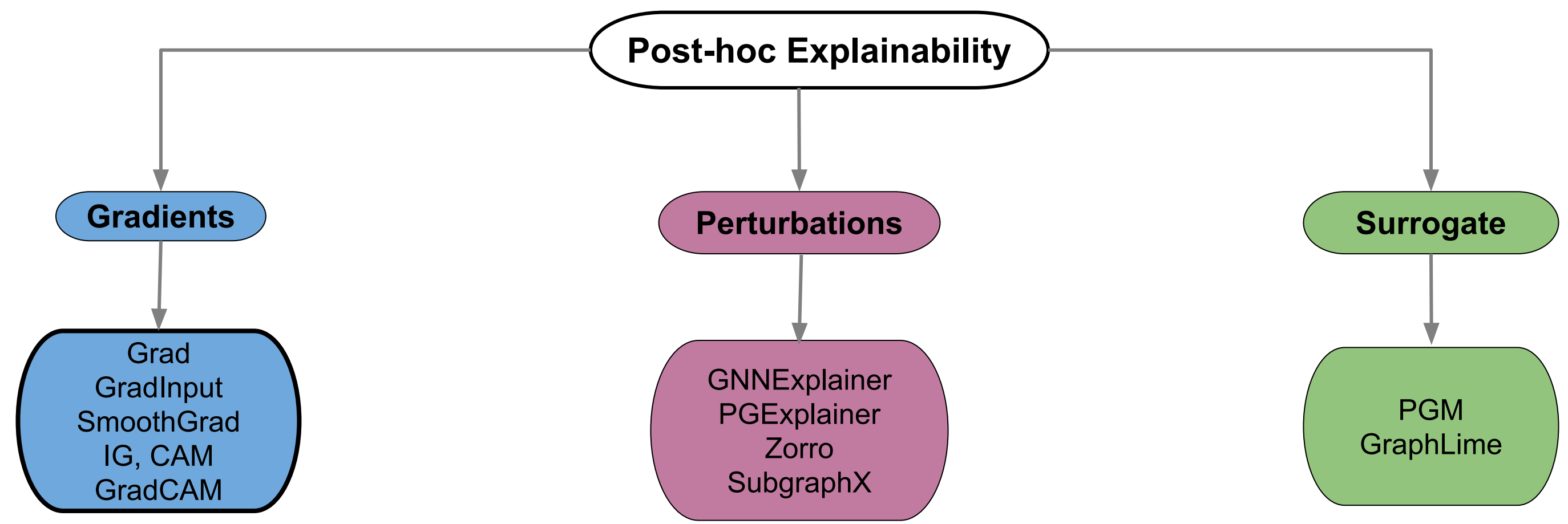


Edge Level



GNN Explanation Methods

Given an input graph $G(A,X)$, a trained GNN Φ , and its prediction



gradient of the prediction with respect to the input as the importance score (mask) for input nodes/edges/features

provide a minimal subgraph of the original graph that is deemed to be important of the prediction

sample a local dataset to represent the relationships around the target node. Usually interpretable by design are applied to fit the sampled local dataset.

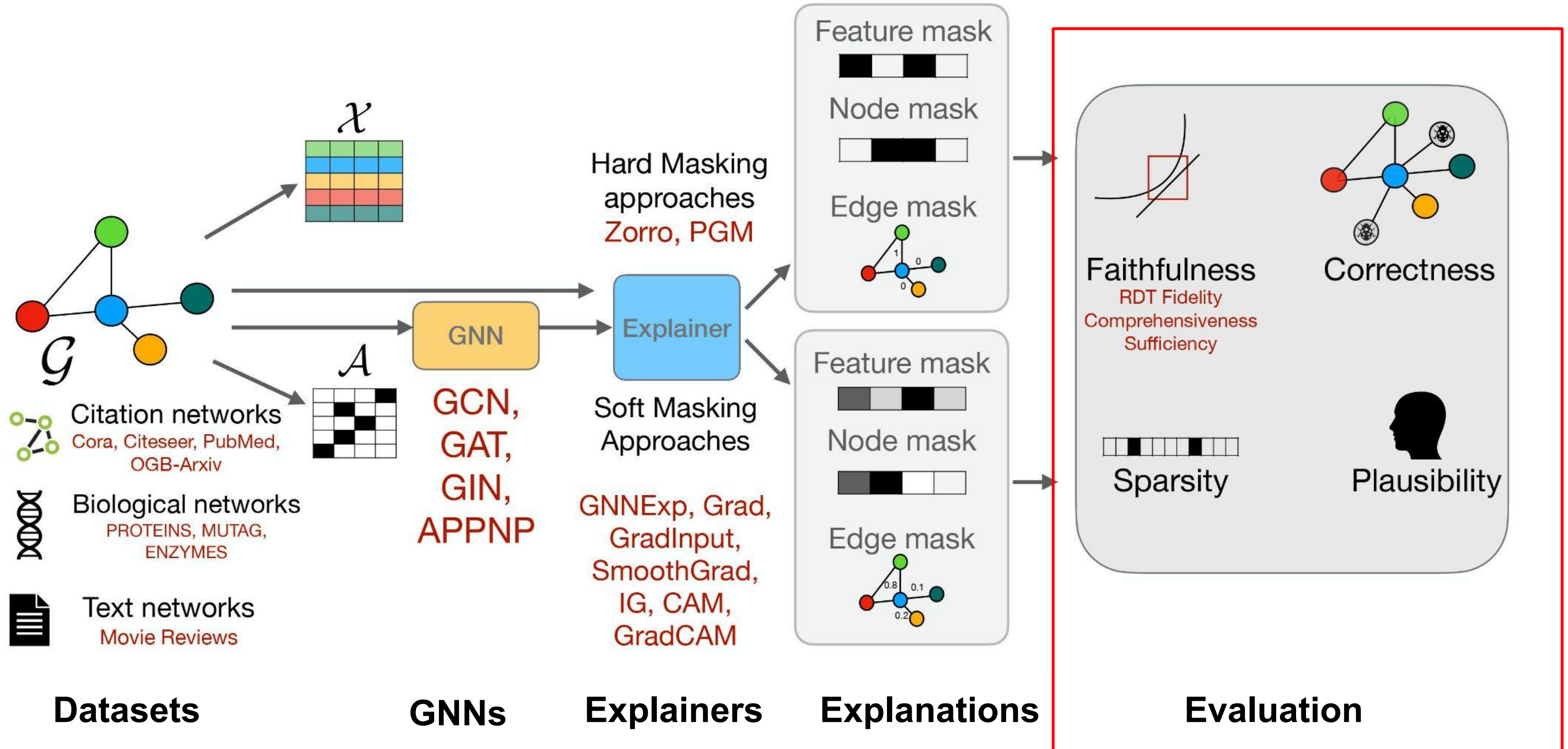
But How to judge if an explanation is good or not?



“a good explanation should be relatively faithful to how the model works, understandable to the receiver, and useful for the receiver’s end-goals.”

- AI Transparency in the Age of LLMs: A Human-Centered Research Roadmap,
(Vera Liao et al.,2023)

Bagel: A Benchmark for Assessing Graph Neural Network Explanations



Faithfulness

How well does the explanation approximate the model's behavior?

The **faithfulness** or **fidelity** of an explanation S can be defined as prediction or accuracy change by removing important nodes/edges/node features.

$$F(S) = \Phi\left(\begin{array}{|c|} \hline \text{yellow} \\ \hline \text{yellow} \\ \hline \text{yellow} \\ \hline \text{yellow} \\ \hline \end{array}\right) - \Phi\left(\begin{array}{|c|} \hline \text{yellow} \\ \hline \text{yellow} \\ \hline \text{yellow} \\ \hline \text{yellow} \\ \hline \end{array} \odot \begin{array}{|c|} \hline 1 \\ \hline 0 \\ \hline 1 \\ \hline 0 \\ \hline \end{array}\right)$$

Drawbacks:

- We do not want to replace the unimportant features with 0, rather its value should not matter.
- Out of distribution features.
- Special case of dense features.

Explainability methods for graph convolutional neural networks (Pope et al., 2019).

Faithfulness

How well does the explanation approximate the model's behavior?

The **faithfulness** or **fidelity** of an explanation S can be defined as prediction or accuracy change by removing important nodes/edges/node features.

RDT-Fidelity

$$F(S) = \Phi\left(\begin{array}{c} \text{yellow} \\ \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{array}\right) - \Phi\left(\begin{array}{c} \text{yellow} \\ \text{yellow} \\ \text{yellow} \\ \text{yellow} \end{array} \odot \begin{array}{c} 1 \\ \text{black} \\ 1 \\ \text{black} \end{array}\right)$$

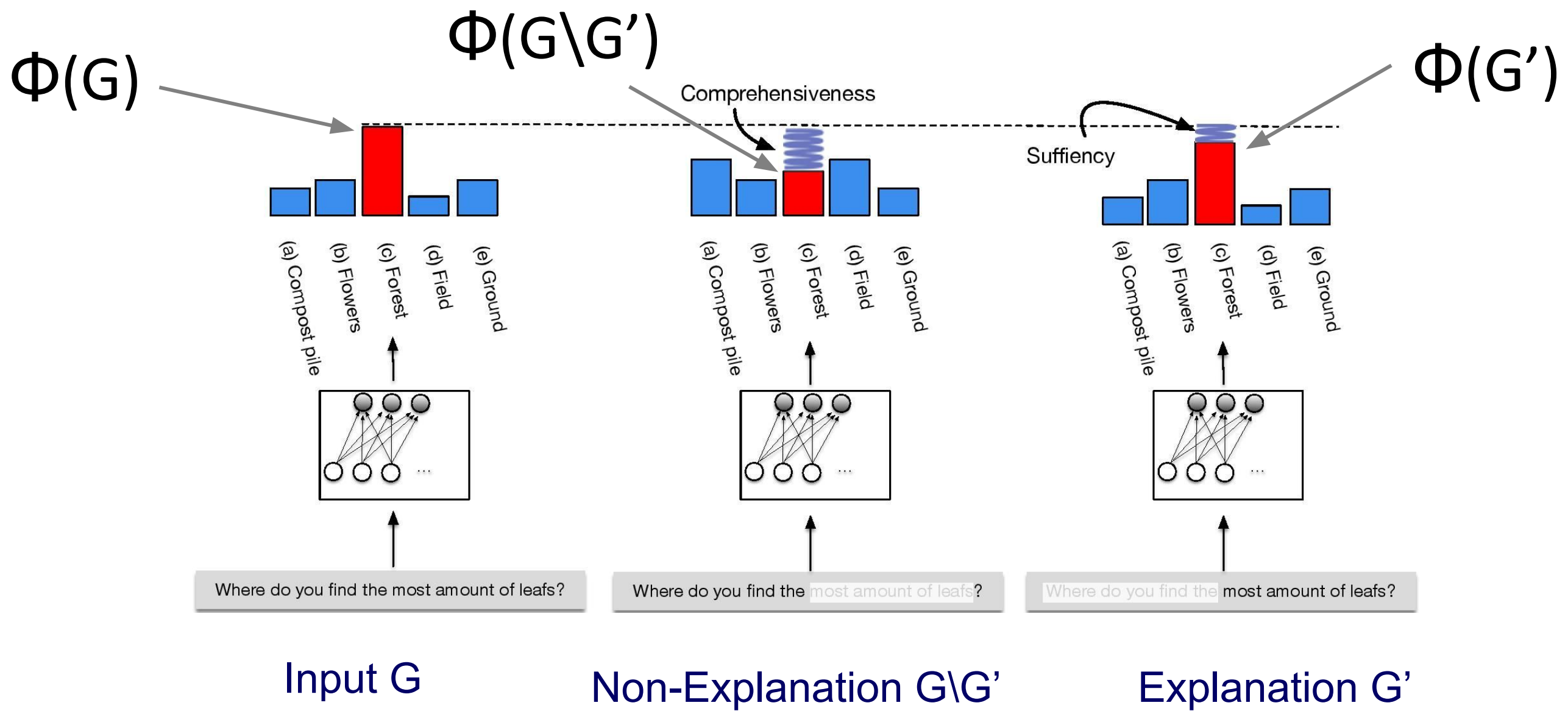
sample from the underline training data distribution

Key Idea: Perturb the unimportant features.

Special case of dense features.

Zorro: Valid, Sparse, and Stable Explanations in Graph Neural Networks (Funke et al., 2022)

Faithfulness



Sufficiency = $\Phi(G) - \Phi(G')$, if the extracted nodes/edges are sufficient to come up the original prediction.

Comprehensiveness = $\Phi(G) - \Phi(G \setminus G')$, if all nodes/edges in the graph needed to make a prediction were selected.

ERASER: DeYoung et al. 2020

But the full input is also a faithful explanation. Are the explanations non-trivial?



Vs.



Sparsity for hard masks (binary explanations) = Selection size/total

Drawbacks:

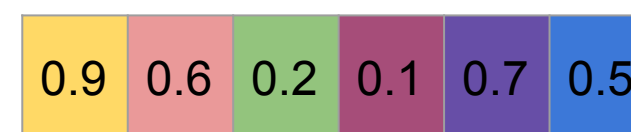
- does not work for soft masks.

Explainability methods for graph convolutional neural networks (Pope et al., 2019).

But the full input is also a faithful explanation. Are the explanations non-trivial?



Vs.

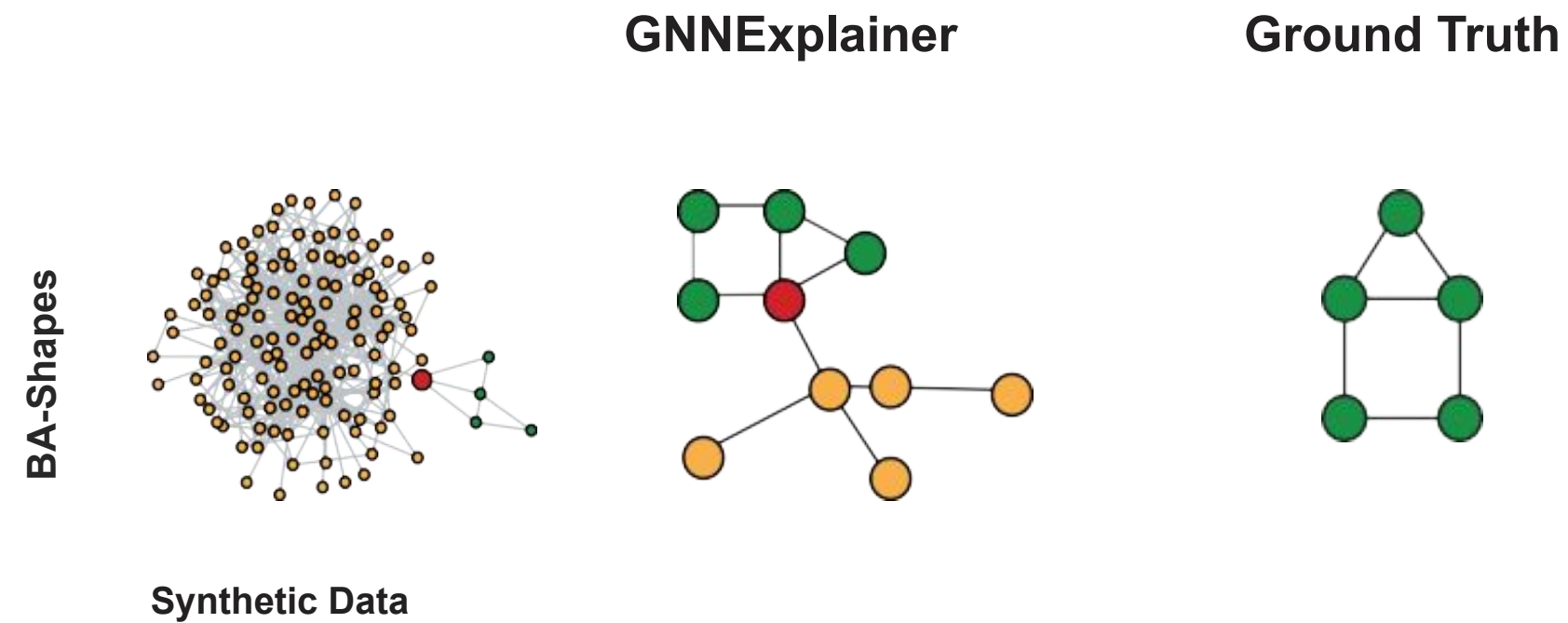


Entropy based Sparsity

$$H(p) = - \sum_x p(x) \log p(x)$$

Explanation Accuracy

How well does the explanation agree with the ground truth?

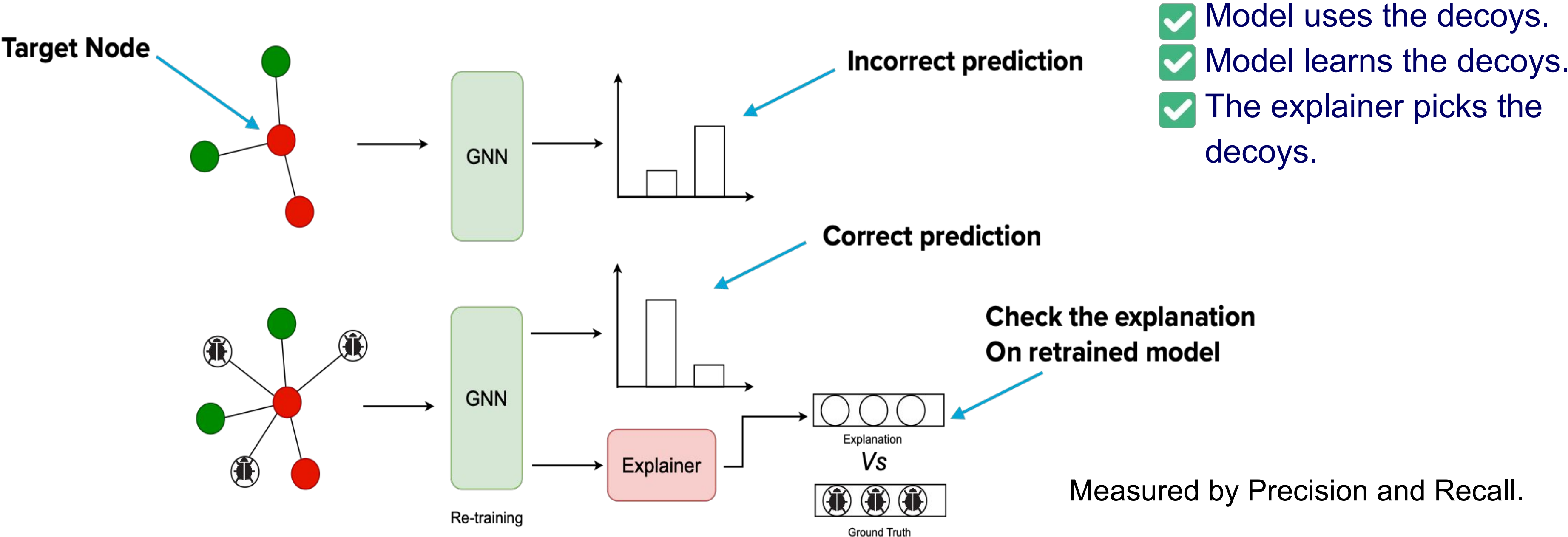


Drawbacks:

- We do not know if the model really used the ground truth for prediction?
- Ground truth is not available for most of the graph datasets.

Correctness

Introduce correlations (decoys) in the training data which can change the decision on a node/graph. Then check if explanation discovers the added correlations.

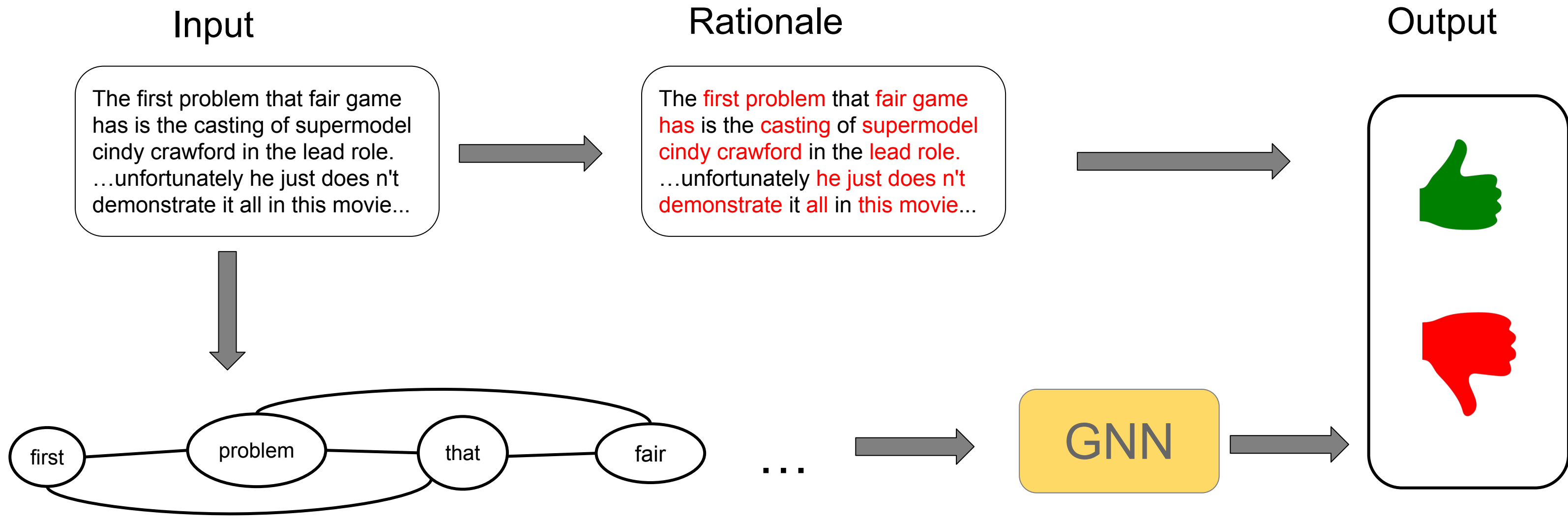


Plausibility

An explanation is considered plausible if it is coherent with human reasoning and understanding.

Graph neural network explanations often disregard plausibility.

Dataset for Plausibility



ERASER: DeYoung et al. 2020

Plausibility

Human Rationales

The **first problem** that **fair game has** is the **casting** of **supermodel cindy crawford** in the **lead role**. not that cindy does that bad... sure william is n't a bad actor. unfortunately **he just does n't demonstrate it all in this movie...**

GNNExp

The first problem that fair game has is the casting of supermodel cindy crawford in the lead role. not that cindy does that **bad**... sure william is n't a **bad actor**. unfortunately he just does n't demonstrate it all in this movie...

Grad

The first problem that fair game has is the casting of supermodel cindy crawford in the lead role. not that cindy does that bad... sure william is **n't a bad actor**. unfortunately **he just does n't demonstrate it all in this movie...**

CAM

The first problem that **fair game has** is the **casting** of **supermodel cindy crawford** in the **lead role**. not that **cindy does** that **bad**... sure william is **n't a bad** actor. unfortunately **he just does n't demonstrate it all in this movie...**

Explainer

Vs.



Measured by

- Token level F1.
- AUPRC(only for soft masks).

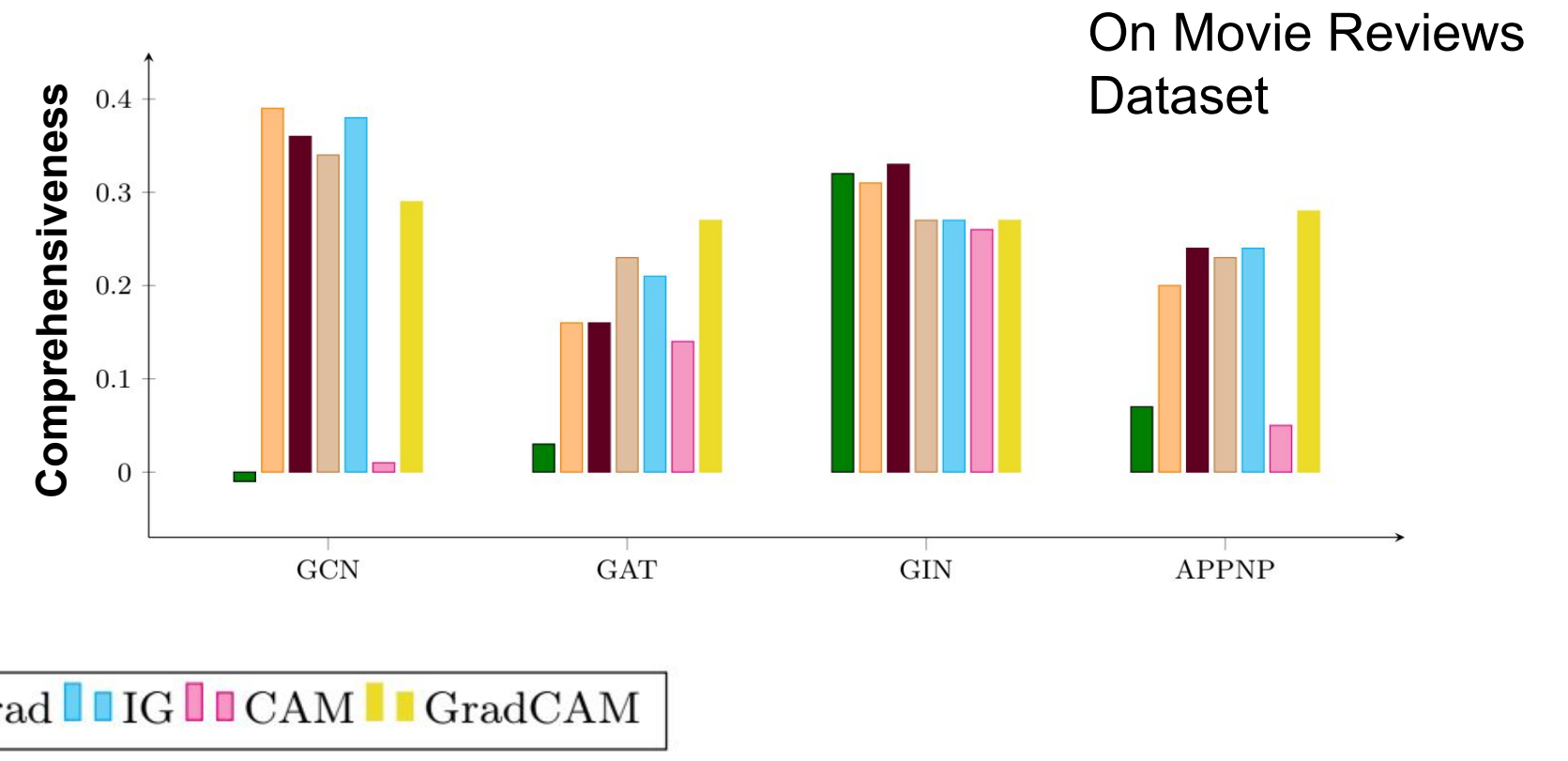
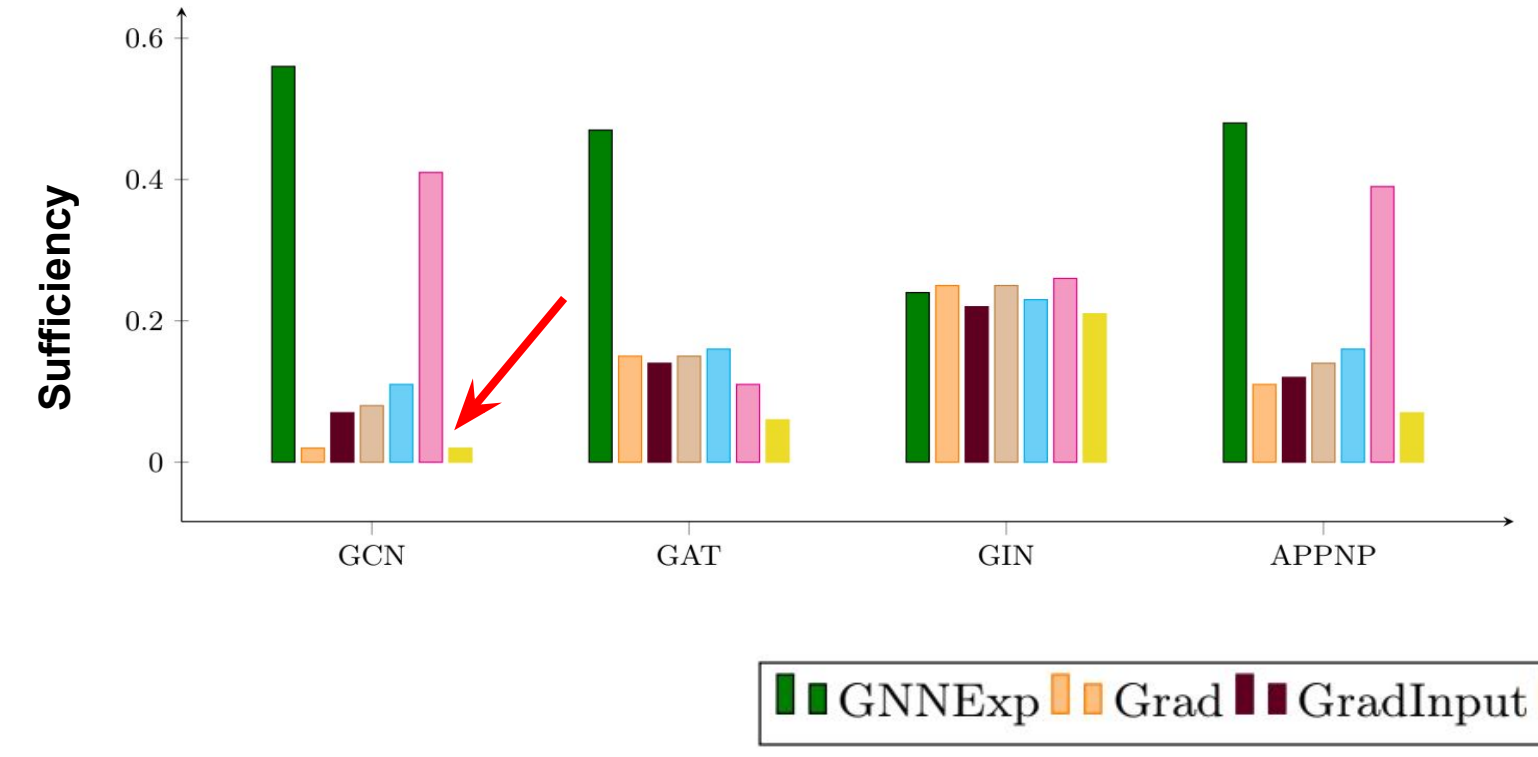
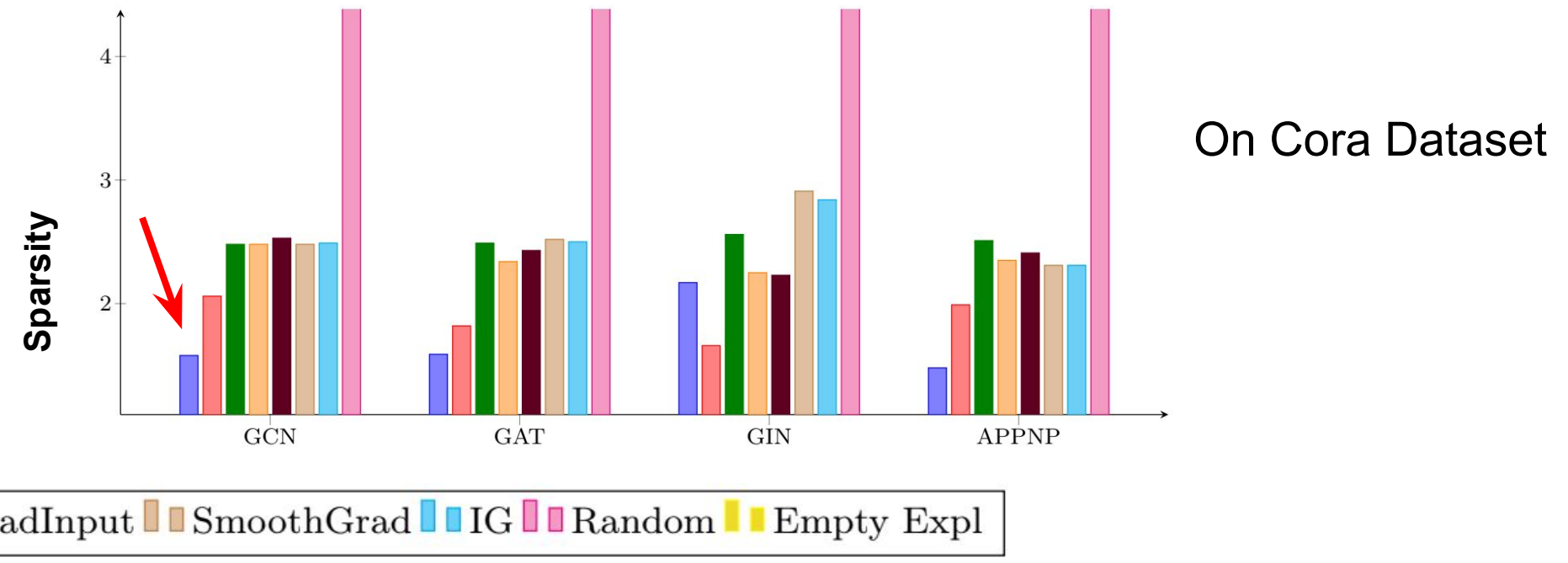
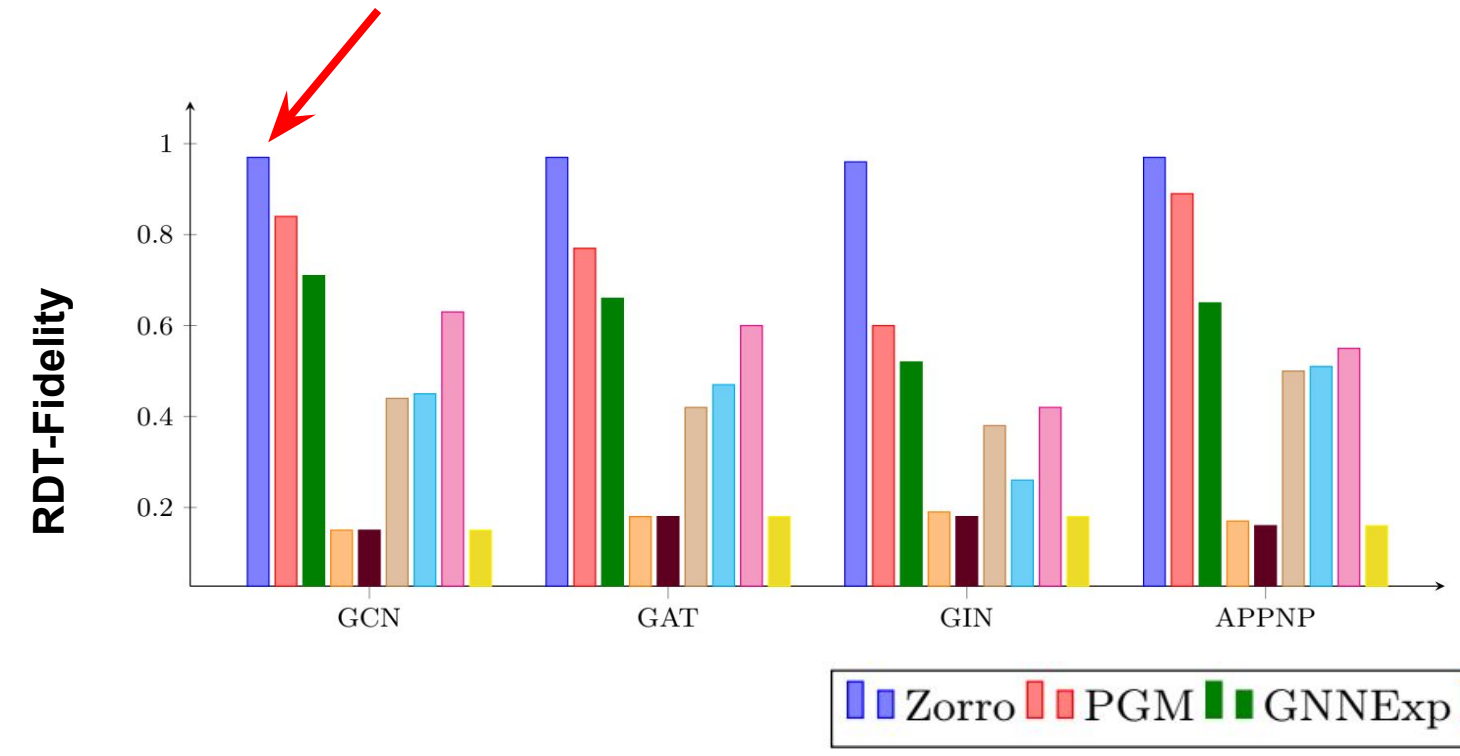
Results

We did ton of experiments in our Bagel benchmark and can be extended easily.

Task	Dataset	Metric				
		<i>Faithfulness</i> RDT- Fidelity	<i>Suff. &</i> Comp.	Sparsity	Correctness	Plausibility
Node	CORA	✓	✗	✓	✓	✗
	CITeseer	✓	✗	✓	✓	✗
	PUBMED	✓	✗	✓	✗	✗
	OGBN-ARXIV	✓	✗	✓	✗	✗
	Synthetic	✓	✗	✓	✓	✓
Graph	Movie Reviews	✗	✓	✗	✗	✓
	MUTAG	✓	✓	✓	✗	✗
	PROTEINS	✓	✓	✓	✗	✗
	ENZYMES	✓	✓	✓	✗	✗

9 explainers
9 Datasets
4 GNNs

Results



Correctness

Mask	Methods	CORA															
		GCN				GAT				GIN				APPNP			
		P@k	R@k	F1	S	P@k	R@k	F1	S	P@k	R@k	F1	S	P@k	R@k	F1	S
Hard	Zorro	0.19	0.80	0.30	45	0.25	0.83	0.37	40	0.26	0.45	0.27	33	0.22	0.79	0.33	38
	PGM	0.11	0.22	0.15	20	0.18	0.36	0.24	20	0.18	0.36	0.25	20	0.19	0.38	0.25	20
	GNNExp	0.42	0.84	0.56	20	0.44	0.88	0.59	20	0.50	1.00	0.67	20	0.34	0.67	0.58	20
	Grad	0.23	0.46	0.31	20	0.29	0.58	0.39	20	0.30	0.60	0.40	20	0.33	0.67	0.45	20
Soft	GradInput	0.16	0.32	0.21	20	0.28	0.56	0.34	20	0.30	0.60	0.40	20	0.28	0.56	0.38	20
	SmoothGrad	0.12	0.25	0.16	20	0.24	0.48	0.32	20	0.50	1.00	0.67	20	0.22	0.43	0.29	20
	IG	0.16	0.32	0.22	20	0.24	0.49	0.33	20	0.50	1.00	0.67	20	0.28	0.55	0.37	20

We added 10 decoys and take k=20.

GNNExp is the best in detecting the injected decoys.

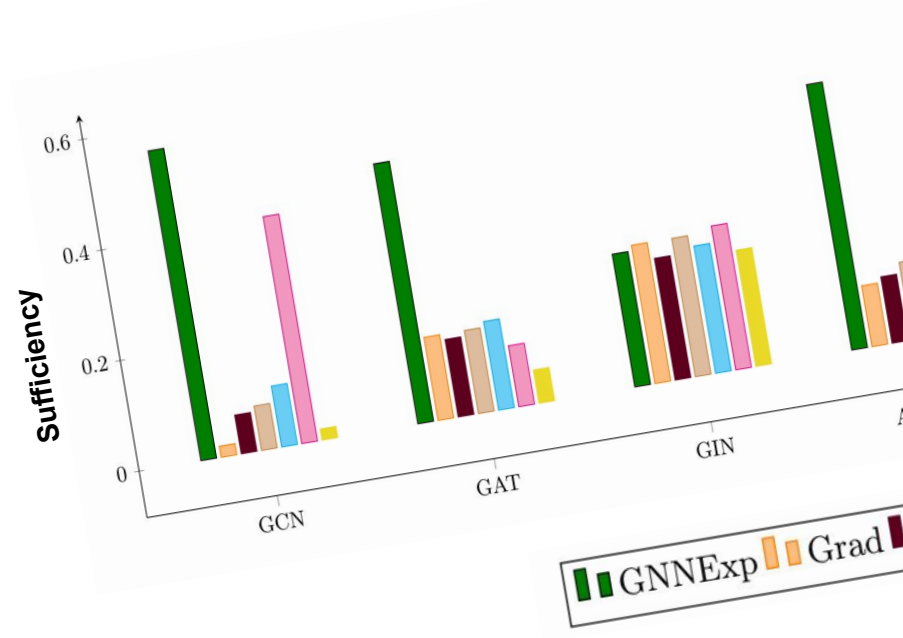
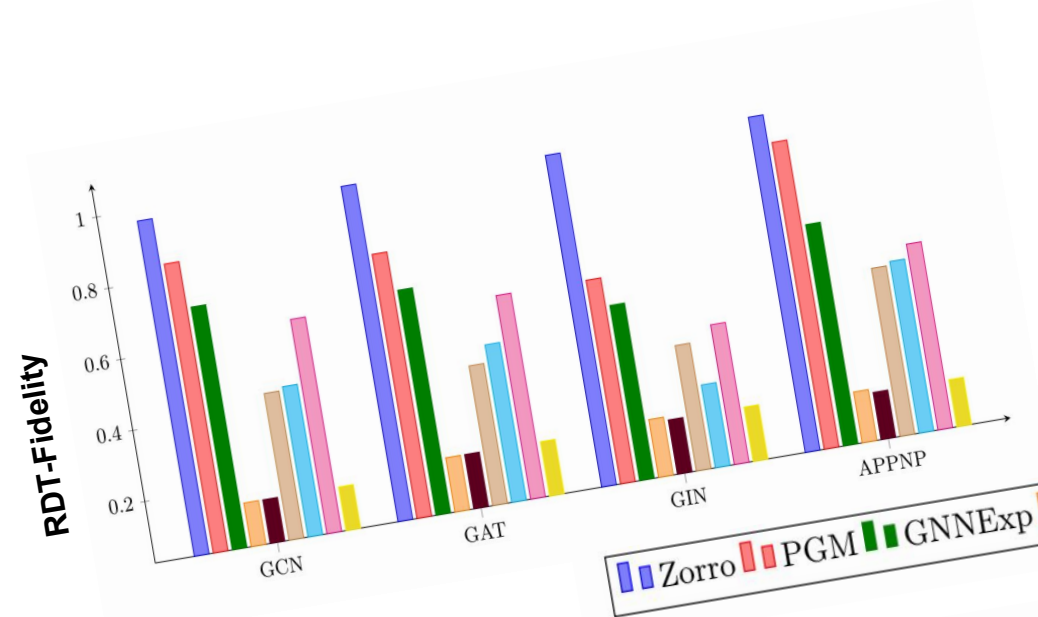
Results

Human Rationales	The first problem that fair game has is the casting of supermodel cindy crawford in the lead role . not that cindy does that bad... sure william is n't a bad actor. unfortunately he just does n't demonstrate it all in this movie...
GNNExp	The first problem that fair game has is the casting of supermodel cindy crawford in the lead role. not that cindy does that bad ... sure william is n't a bad actor . unfortunately he just does n't demonstrate it all in this movie...
Grad	The first problem that fair game has is the casting of supermodel cindy crawford in the lead role. not that cindy does that bad... sure william is n't a bad actor . unfortunately he just does n't demonstrate it all in this movie...
CAM	The first problem that fair game has is the casting of supermodel cindy crawford in the lead role . not that cindy does that bad ... sure william is n't a bad actor. unfortunately he just does n't demonstrate it all in this movie...

Mask	Methods	GCN		
		auprc	F1	S
Hard	PGM	—	0.42	25
	GNNExp	0.46	0.54	168
Soft	Grad	0.44	0.52	265
	GradInput	0.39	0.51	221
	SmoothGrad	0.40	0.52	219
	IG	0.37	0.49	225
	CAM	0.54	0.61	224
	GradCAM	<u>0.67</u>	0.34	175

Plausibility

Results



Mask Methods	GCN		GAT		CORRA		GIN		APPNP			
	P@k	R@k	F1	S	P@k	R@k	F1	S	P@k	R@k	F1	S
Zorro	0.37	0.40	0.26	0.45	0.27	0.33	0.22	0.79	0.33	0.38	0.25	20
PGM	0.24	0.20	0.18	0.36	0.25	0.20	0.19	0.38	0.25	0.20	0.25	20
GNNExp	0.59	0.20	0.50	1.00	0.67	0.20	0.34	0.67	0.58	0.20	0.20	20
Zorro	0.39	0.20	0.30	0.60	0.40	0.20	0.33	0.67	0.45	0.20	0.45	20
PGM	0.34	0.20	0.30	0.60	0.40	0.20	0.28	0.56	0.38	0.20	0.38	20
GNNExp	0.32	0.20	0.50	1.00	0.67	0.20	0.22	0.43	0.29	0.20	0.29	20
Zorro							0.28	0.55	0.37	0.20	0.37	20

A cartoon by Bill Abbott showing two people sitting on the ground, looking at a dead animal. The scene is set in a desolate, rocky landscape. The cartoon is signed 'Bill Abbott' in the bottom right corner.

“A truly lively debate, yet no clear winner”
- Bill Abbott

the casting of supermodel cindy crawford in the lead role. not n't a bad actor. unfortunately he just does n't demonstrate it

casting of supermodel cindy crawford in the lead role. not n't a bad actor. unfortunately he just does n't demonstrate it

asting of supermodel cindy crawford in the lead role. not a bad actor. unfortunately he just does n't demonstrate it

sting of supermodel cindy crawford in the lead role. not bad actor. unfortunately he just does n't demonstrate it

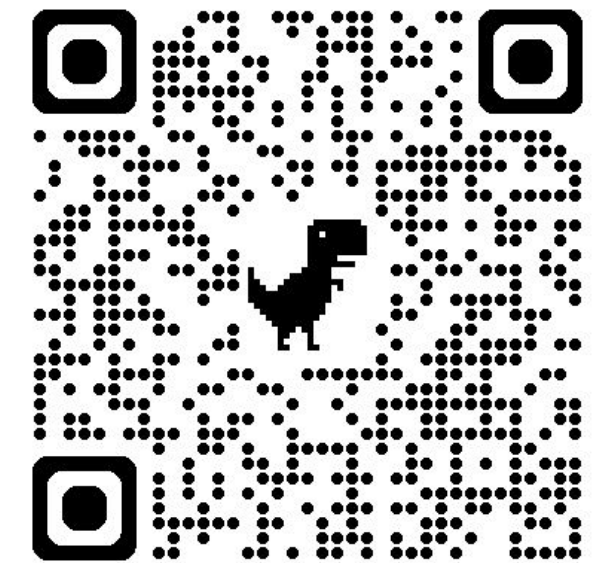
Conclusion

- Explainability in Graph
- Four dimensions of measuring the goodness of the explanation
 - Faithfulness
 - Sparsity
 - Correctness
 - Plausibility
- New dataset for Plausibility.
- Study different explainers.
- Lot more to explore



Thank You

 rathee@l3s.de



Please scan for the
Github link