

---

# Graphs Counterfactual Explainability: A Comprehensive Landscape

Tutorial @ AAAI 2024

M.A. Prado-Romero - B. Prenkaj - G. Stilo





Mario A. Prado-Romero  
PhD. Candidate  
*Gran Sasso Science Institute*



Bardh Prenkaj  
PostDoc  
*Sapienza University of Rome*



Giovanni Stilo  
Associate Prof.  
*University of L'Aquila*



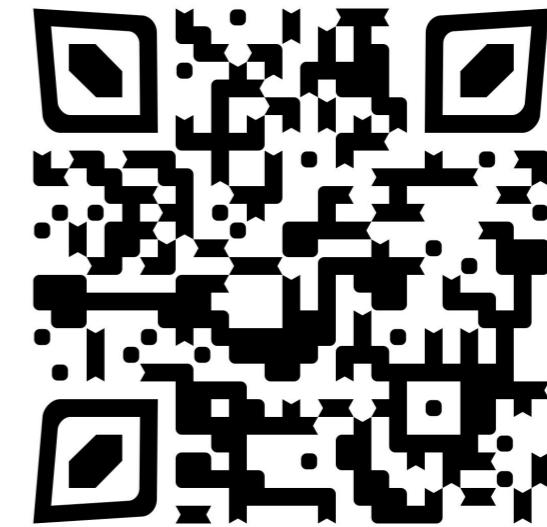
We are with the *Artificial Intelligence & Information Mining* (**aiim** - pronounced as i'm /aɪm/, and aim /eɪm/) a collective of *Individuals* (/aɪm/) who share a common *Interest* (/eɪm/) in Artificial Intelligence, Data Mining, and Machine Learning

# ROADMAP

---

- **Part I: Graphs fundamentals (20 mins) [STILO]**
  - Wide-spread adoption of GNNs in graph prediction problems
  - How do GNNs work?
  - Applications of different types of GNNs
- **Part II: eXplainable AI (25 mins) [PRENKAJ]**
  - Issues of black-box models and the importance of interpretability
  - What is a factual explanation?
  - (Briefly) Revisiting GNNExplainer and GraphLIME
- **Part III: Counterfactual Explanations in Graphs (60 mins) [PRADO]**
  - What is a graph counterfactual explanation (GCE)?
  - GCE taxonomy description and method classification
    - Instance-level explainers
      - Search-based
      - Heuristic-based
      - Learning-based
    - Model-level explainers
  - Benchmarking datasets and evaluation metrics (pro et contra)

based on ACM survey



<https://dl.acm.org/doi/abs/10.1145/3618105>



PART I

# GRAPHS FUNDAMENTALS AND THEIR NEURAL NETWORK

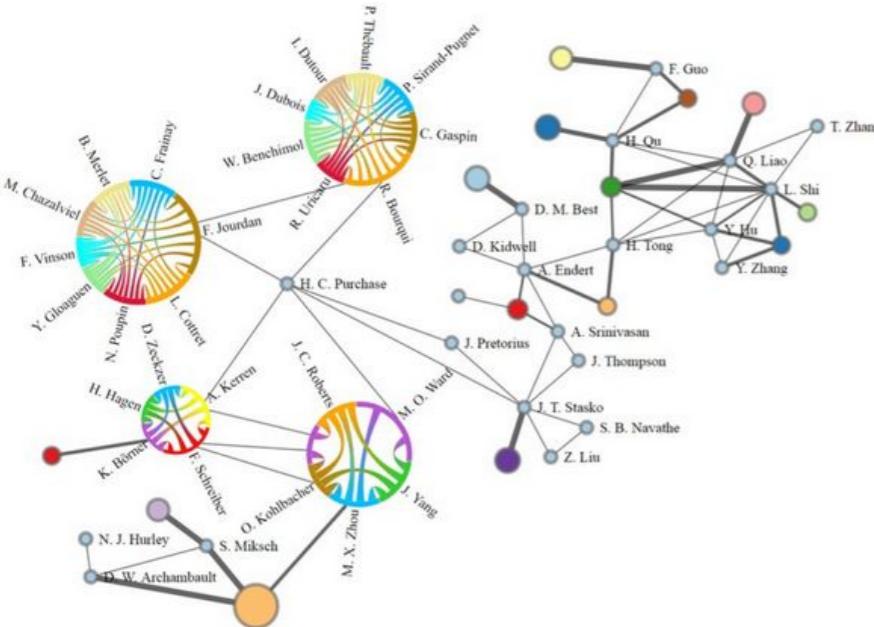
by Giovanni Stilo

images based on Understanding Deep Learning - book by Simon J.D. Prince

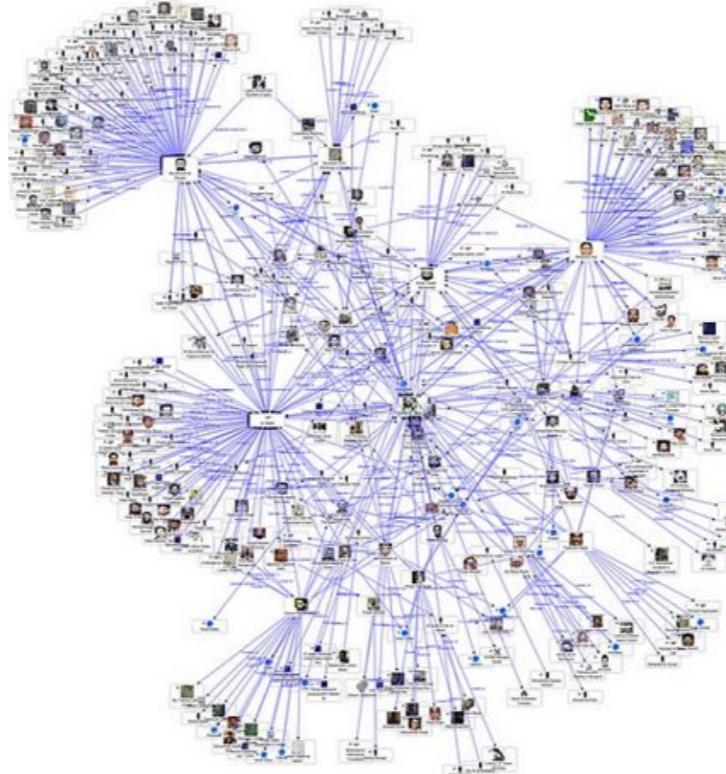


# WHAT IS A GRAPH?

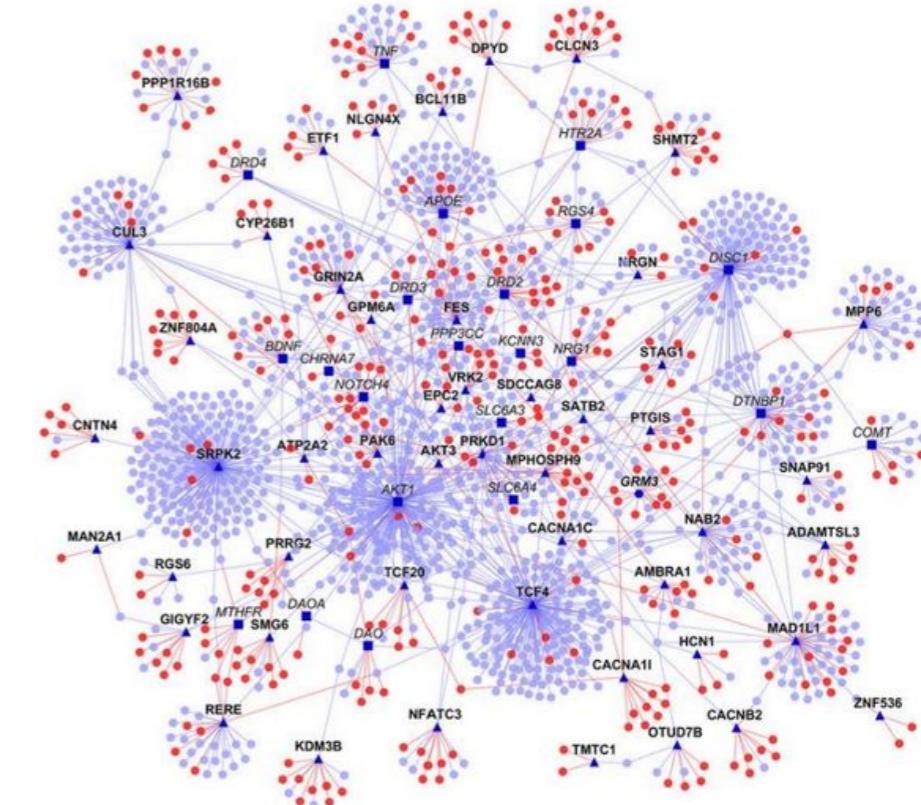
---



CO-AUTORSHIP NETWORKS



SOCIAL NETWORKS



PROTEIN INTERACTIONS

# CHALLENGES WITH GRAPHS

---

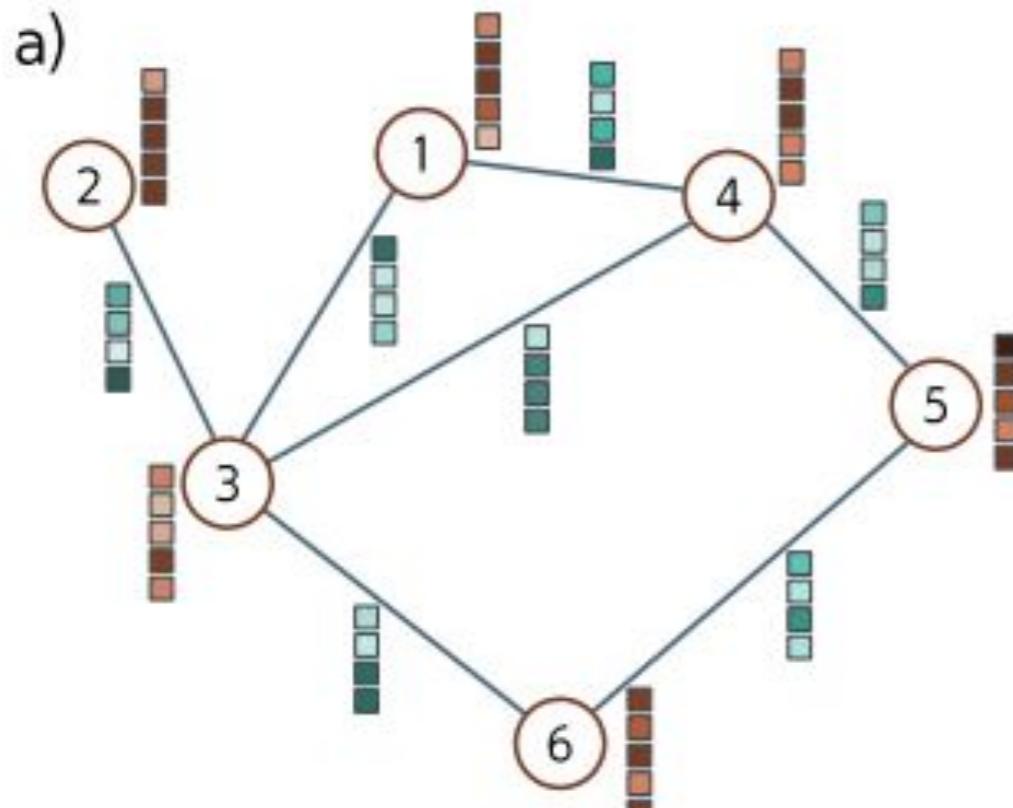
There are **three main challenges** associated with processing graphs:

- **Variable topology:** hard to design an NN that is sufficiently expressive and can cope with this variation
- **Huge graphs:** we can have millions of nodes and billions of edges (see Twitter)
- **Single monolithic graph:** the usual protocol of training with many data examples and testing with new data is not always appropriate or possible



# GRAPH REPRESENTATION

**adjacency matrix**,  $A$  is  $N \times N$  ;  
**node embeddings**,  $X$  is  $D \times N$  ;  
**edges embeddings**,  $E$  is  $D_E \times N$



b)

Adjacency  
matrix,  $A$   
 $N \times N$

	1	2	3	4	5	6
1	□	□	□	■	□	□
2	□	□	■	■	□	□
3	■	■	□	■	■	■
4	■	■	■	□	■	■
5	■	■	■	■	□	■
6	■	■	■	■	■	□

c)

Node  
data,  $X$   
 $D \times N$

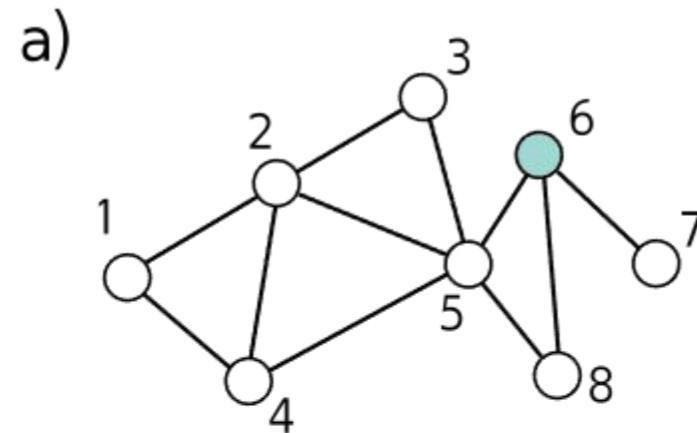
1	2	3	4	5	6
■	■	■	■	■	■
■	■	■	■	■	■
■	■	■	■	■	■
■	■	■	■	■	■
■	■	■	■	■	■

d)

Edge  
data,  $E$   
 $D_E \times E$

1	1	2	3	3	4	5
3	4	3	4	6	5	6
■	■	■	■	■	■	■
■	■	■	■	■	■	■
■	■	■	■	■	■	■
■	■	■	■	■	■	■

# ADJACENCY MATRIX PROPERTIES



b)

$$\mathbf{A} = \begin{bmatrix} 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 \end{bmatrix}$$

c)

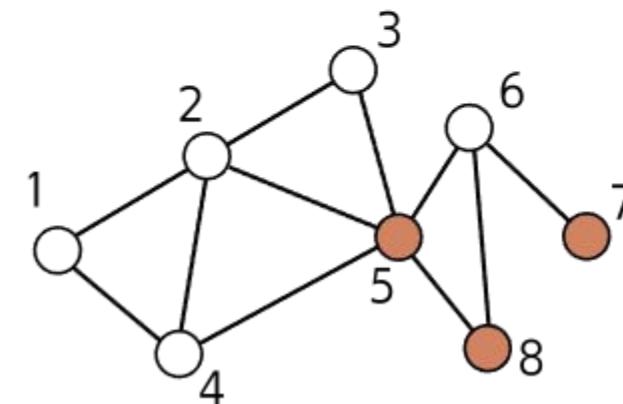
$$\mathbf{A}^2 = \begin{bmatrix} 2 & 1 & 1 & 1 & 2 & 0 & 0 & 0 \\ 1 & 4 & 1 & 2 & 2 & 1 & 0 & 1 \\ 1 & 1 & 2 & 2 & 1 & 1 & 0 & 1 \\ 1 & 2 & 2 & 3 & 1 & 1 & 0 & 1 \\ 2 & 2 & 1 & 1 & 5 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 3 & 0 & 1 \\ 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 & 1 & 1 & 2 \end{bmatrix}$$

d)

$$\mathbf{x} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

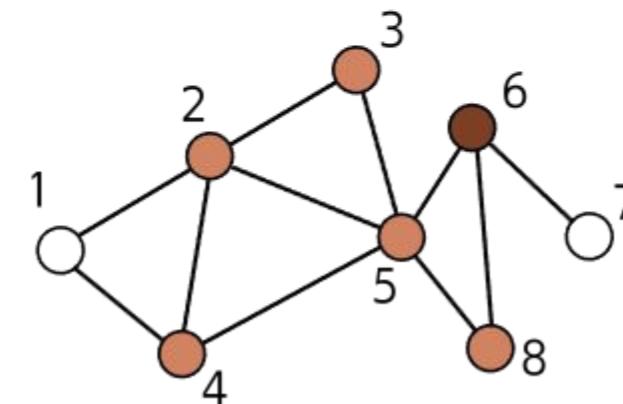
e)

$$\mathbf{Ax} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 0 \\ 1 \end{bmatrix}$$



f)

$$\mathbf{A}^2\mathbf{x} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 1 \\ 3 \\ 0 \\ 1 \end{bmatrix}$$

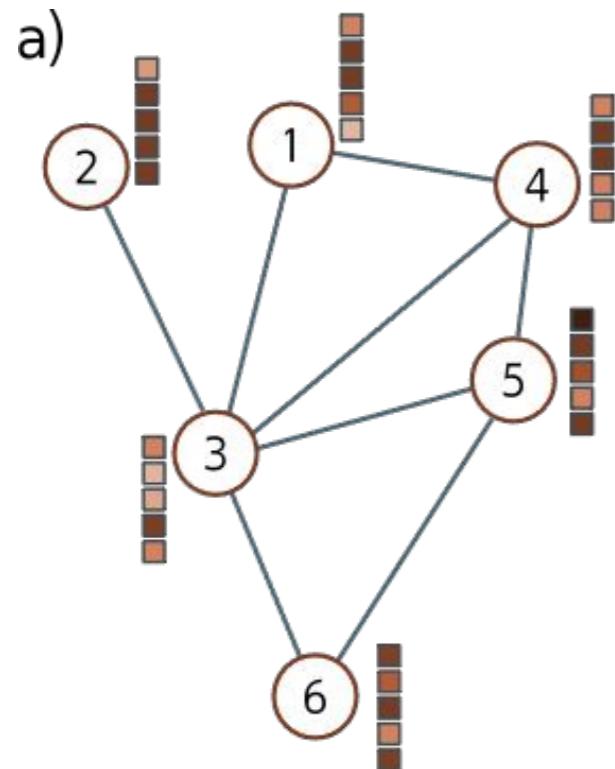


$\mathbf{A}^l$  contains the number of unique walks of length  $l$  from node  $m$  to node  $n$

# NODES PERMUTATION

Node indexing in graphs is arbitrary

Permuting the node **indices** results in a **permutation** of the **columns** of the node data matrix  $X$  and a permutation of both the **rows** and **columns** of the adjacency matrix  $A$ .

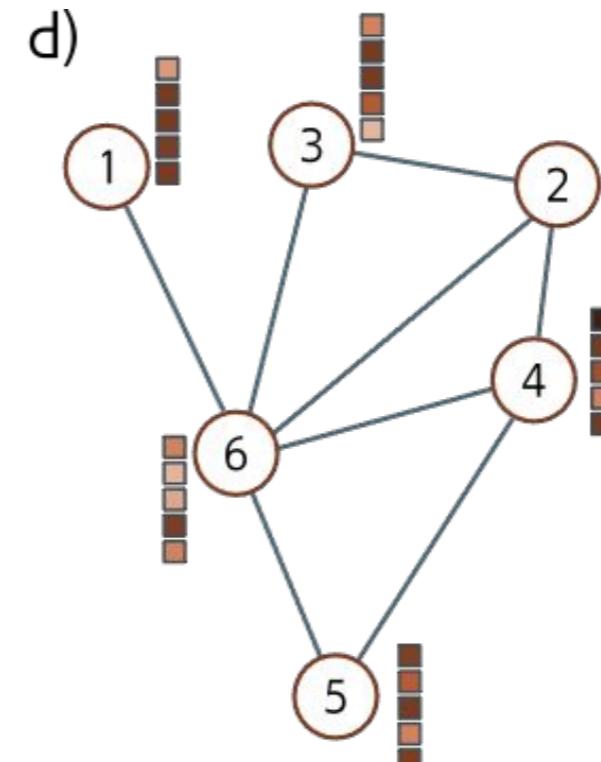


b) Adjacency  $A$

	1	2	3	4	5	6
1	0	0	1	1	0	0
2	0	0	1	0	0	0
3	1	1	0	1	1	1
4	1	0	0	0	1	1
5	0	0	1	1	0	1
6	0	0	0	1	1	0

c) Node data,  $X$

	1	2	3	4	5	6
1	1	2	3	4	5	6
2	2	3	4	5	6	1
3	3	4	5	6	1	2
4	4	5	6	1	2	3
5	5	6	1	2	3	4
6	6	1	2	3	4	5



e) Adjacency  $A$

	1	2	3	4	5	6
1	0	0	1	1	0	0
2	0	0	1	0	0	0
3	1	1	0	1	1	1
4	1	0	0	0	1	1
5	0	0	1	1	0	1
6	0	0	0	1	1	0

f) Node data,  $X$

	1	2	3	4	5	6
1	1	2	3	4	5	6
2	2	3	4	5	6	1
3	3	4	5	6	1	2
4	4	5	6	1	2	3
5	5	6	1	2	3	4
6	6	1	2	3	4	5

$$X' = X \mathbf{P}, \quad A' = \mathbf{P}^T A \mathbf{P}$$

# GRAPH LEARNING

---

We want to learn a (dense) representation  $\mathbf{H}$  of the graph usable for different downstream tasks

A **graph neural network** is a model that takes:

- the node embeddings  $\mathbf{X}$  and the adjacency matrix  $\mathbf{A}$  as inputs and passes them through a series of  $k$  layers.
- the node embeddings are updated at each layer to create intermediate “hidden” representations  $\mathbf{h}$  before finally computing output embeddings  $\mathbf{h}_K$ .

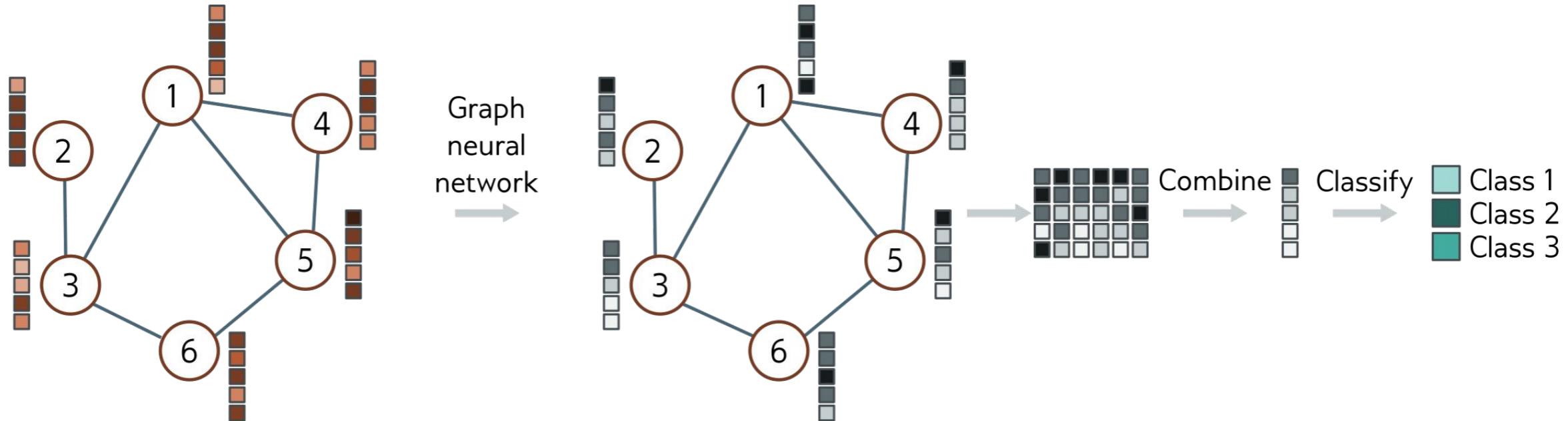


# GRAPH CLASSIFICATION TASKS

For example, we might want to **predict**:

- the **temperature** at which a **molecule** becomes **liquid** (a regression task);
- whether a **molecule** is **poisonous** to human beings or not (a classification task).

For graph-level tasks, the output **node embeddings are combined** (e.g., by averaging), and the resulting vector is **mapped** via a linear transformation or neural network to a **fixed-size vector**.



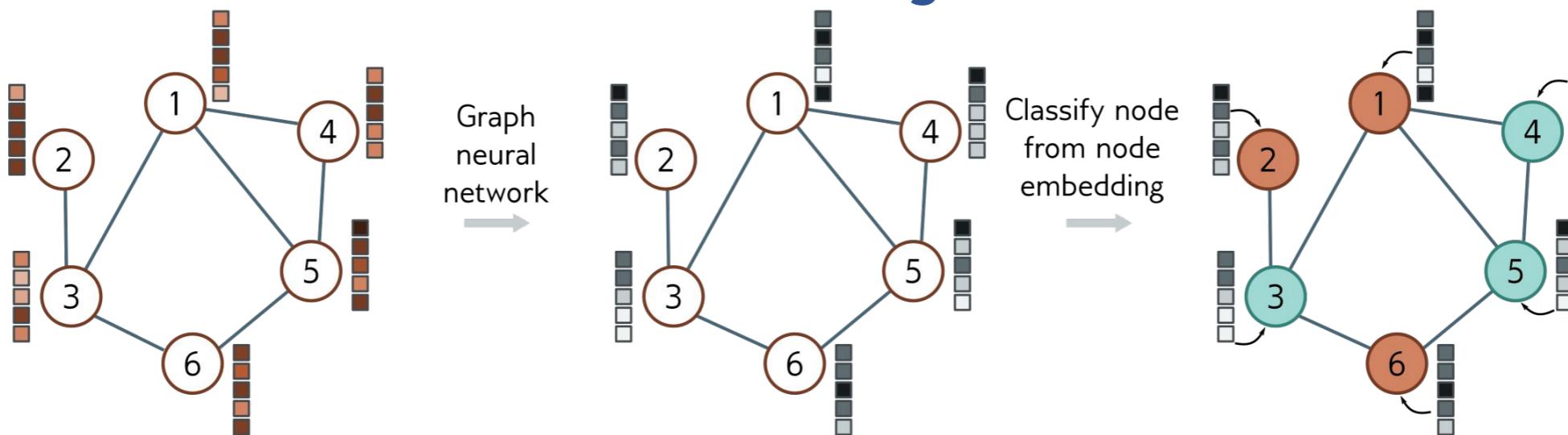
$$Pr(y = 1 | \mathbf{X}, \mathbf{A}) = sig \left( \beta_k + \boldsymbol{\omega}_k \mathbf{H}_k \frac{1}{N} \right)$$

# NODE CLASSIFICATION TASKS

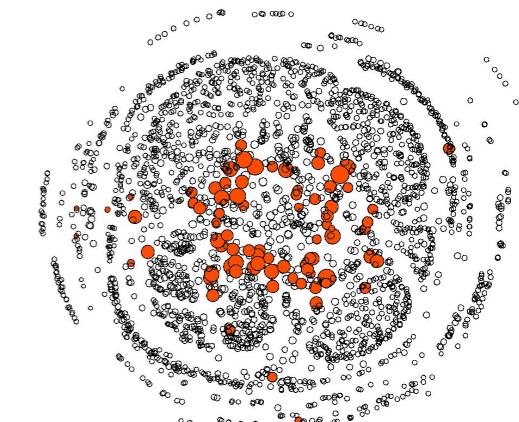
For example, in an *PPI* network we might want to **predict**:

- the probability that a **given node** might be **attacked/being part** of a certain **disease** (classification) as it is shown for COVID19 (**red**) - PPI on the right.

The network assigns **one or more label** (classification) or values (regression) to **each node** of the graph, **using** both the **graph structure** and learned **node embeddings**.



$$Pr(y^{(n)} = 1 | \mathbf{X}, \mathbf{A}) = sig \left( \beta_k^{(n)} + \boldsymbol{\omega}_k^{(n)} \mathbf{H}_k^{(n)} \right)$$

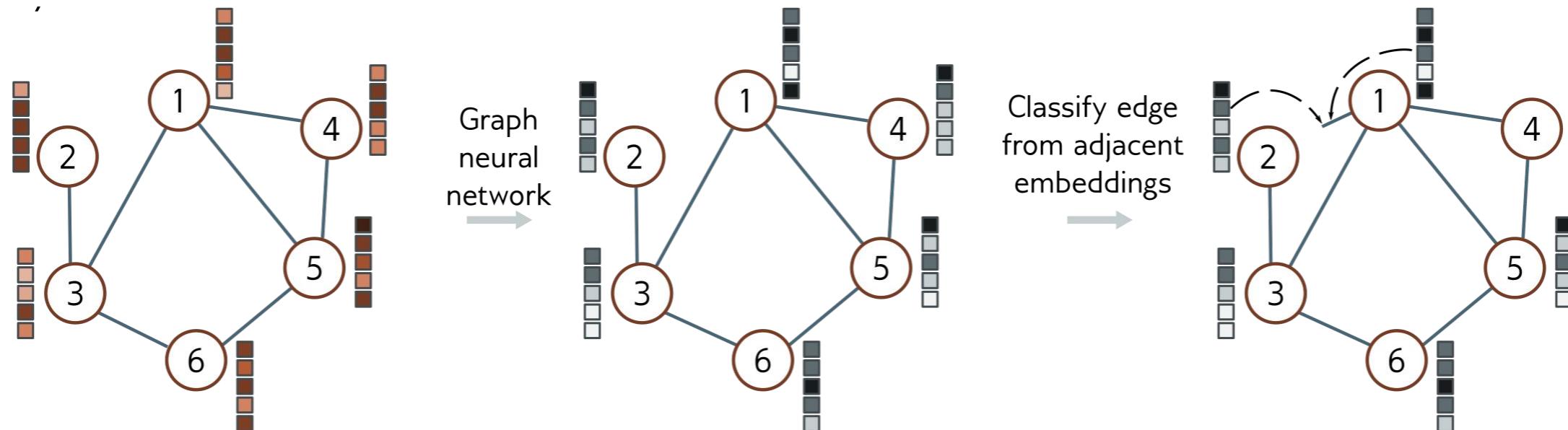


# EDGE CLASSIFICATION TASKS

For example, in the **social network** setting, the network might **predict** whether:

- two people know each other and suggest that they connect if that is the case.

The network assigns **one or more label** (classification) or values (regression) to **each edges** of the graph, using both the graph **structure** and learned **node embeddings**.



$$Pr(y^{(mn)} = 1 \mid \mathbf{X}, \mathbf{A}) = sig \left( \mathbf{H}_k^{(m)T} \cdot \mathbf{H}_k^{(n)} \right)$$

# GNNS

---

**Spatial-based convolutional graph neural networks (GCN)** are convolutional network that update each node embeddings by **aggregating information from nearby nodes** using the original **graph structure**.

Each **layer** of the GCN is a function  $F[\cdot]$  with parameters  $\Phi$  that **takes** the node **embeddings** and **adjacency matrix** and outputs new node embeddings:

$$\mathbf{H}_1 = F[\mathbf{X}, \mathbf{A}, \phi_0]$$

$$\mathbf{H}_2 = F[\mathbf{H}_1, \mathbf{A}, \phi_1]$$

$$\mathbf{H}_3 = F[\mathbf{H}_2, \mathbf{A}, \phi_2]$$

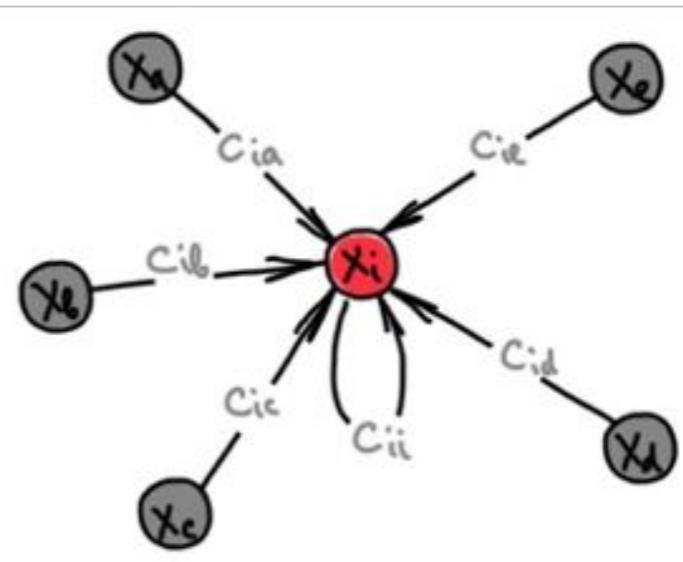
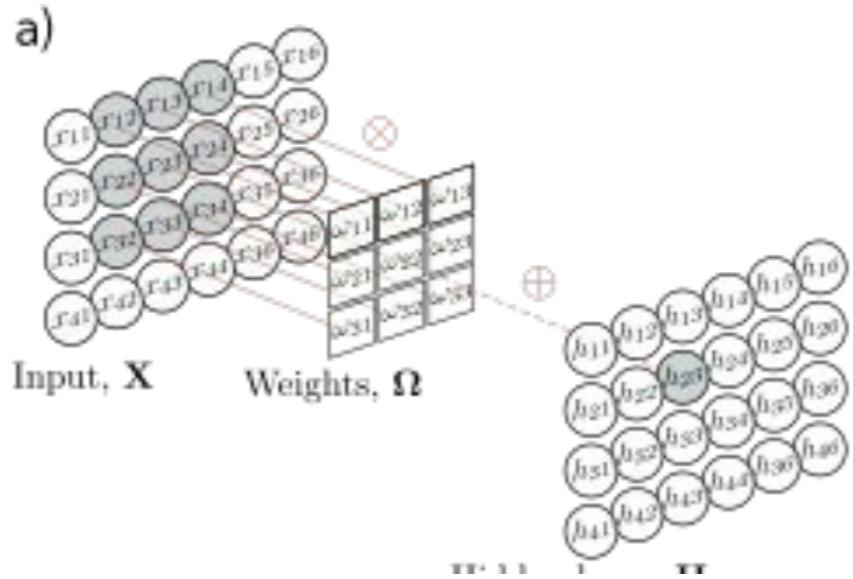
$$\vdots = \vdots$$

$$\mathbf{H}_K = F[\mathbf{H}_{K-1}, \mathbf{A}, \phi_{K-1}]$$



# PARAMETER SHARING

- Likewise CNN, we want the **same parameters** at every **node**: **reducing** the number of **parameters** and **sharing** what the network **learns** at each node **across** the **entire** graph.



- Each neighbor sends a message to the variable of interest, which aggregates these messages to form the update.

In **images**, the neighbors were **pixels** from a **fixed-size square region** around the current position, so the **spatial relationships** at each position **are the same**.

- In a **graph**, each node may have a **different number of neighbors**, and there are **no consistent relationships**.

# GCN (1)

each node at layer  $k$ , we aggregate information from neighboring nodes by e.g. summing their node embeddings:

$$\mathbf{agg}[n, k] = \sum_{m \in \text{ne}[n]} \mathbf{H}_k^{(m)}$$

linear transformation  $\Omega$  to the embedding  $H$  of the current node and to his aggregated value, we add a bias term  $\beta$ , and pass the result through a nonlinear activation function  $a[\cdot]$ , which is applied independently to every member of its vector argument:

$$\mathbf{H}_{k+1}^{(n)} = a \left[ \boldsymbol{\beta}_k + \boldsymbol{\Omega}_k \cdot \mathbf{H}_k^{(n)} + \boldsymbol{\Omega}_k \cdot \mathbf{agg}[n, k] \right]$$

the  $n^{th}$  column of  $\mathbf{A}$  contains ones at the positions of neighbors. If post-multiply the embeddings by  $\mathbf{A}$  the  $n^{th}$  column is  $\mathbf{agg}[n, k]$ :

$$\mathbf{H}_{k+1} = a[\boldsymbol{\beta}_k \mathbf{1}^T + \boldsymbol{\Omega}_k \mathbf{H}_k + \boldsymbol{\Omega}_k \mathbf{H}_k \mathbf{A}] = a[\boldsymbol{\beta}_k \mathbf{1}^T + \boldsymbol{\Omega}_k \mathbf{H}_k (\mathbf{A} + \mathbf{I})]$$



# GCN (2)

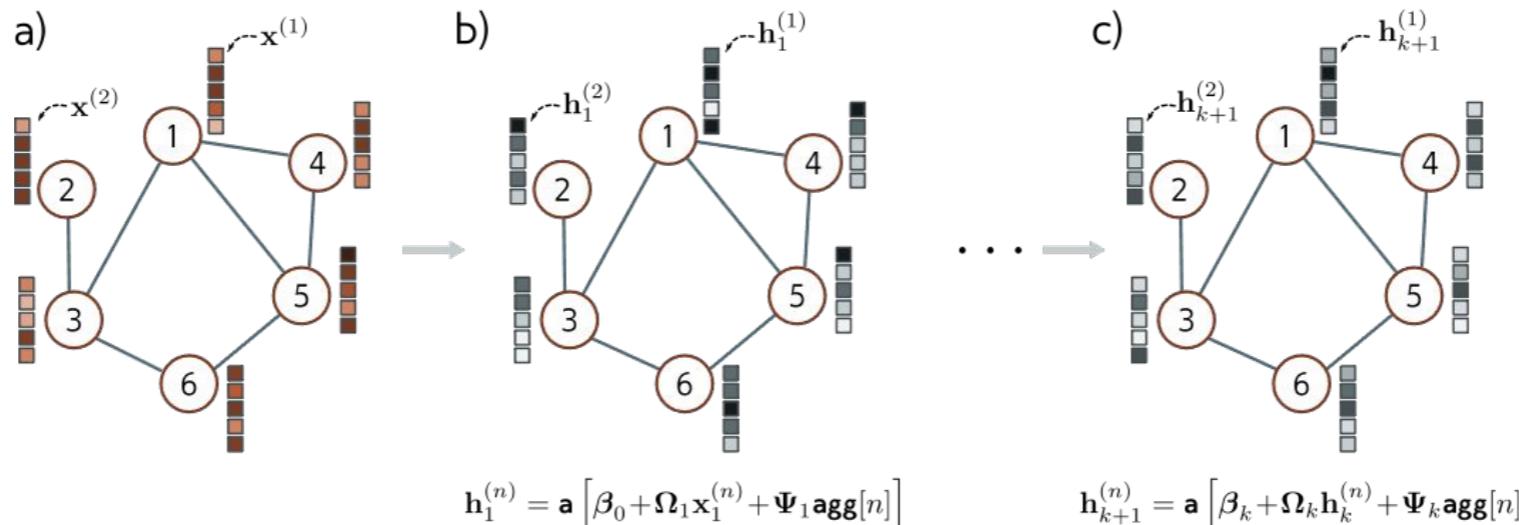
$$\mathbf{H}_{k+1} = a[\boldsymbol{\beta}_k \mathbf{1}^T + \boldsymbol{\Omega}_k \mathbf{H}_k (\mathbf{A} + \mathbf{I})]$$

This layer **satisfies** the **design** considerations:

- it **is equivariant to permutations** of the node **indices**

$$\mathbf{H}_{k+1}\mathbf{P} = \text{F}[\mathbf{H}_k\mathbf{P}, \mathbf{P}^T \mathbf{A} \mathbf{P}, \phi_k] = a[\boldsymbol{\beta}_k \mathbf{1}^T \mathbf{P} + \boldsymbol{\Omega}_k \mathbf{H}_k \mathbf{P} (\mathbf{P}^T \mathbf{A} \mathbf{P} + \mathbf{I})]$$

- can **cope** with **any** number of **neighbors** due to the **agg[n,k]** function;
- **exploits** the **graph structure** to provide a relational inductive bias,
- and **shares parameters** throughout the graph i.e.  **$\boldsymbol{\Omega}$** .



# GRAPH CLASSIFICATION (REV.)

- We want a **neural network**  $f[\mathbf{X}, \mathbf{A}, \Phi]$  that **classifies** (predicts) **molecules** as toxic or harmless.
- The **adjacency matrix**  $\mathbf{A} \in \mathbb{R}^{N \times N}$  derives from the molecular structure.
- The columns of the **node embedding** matrix are **one-hot vectors** indicating which of the 118 **elements of the periodic** table are present.

$$\mathbf{H}_1 = a[\boldsymbol{\beta}_0 \mathbf{1}^T + \boldsymbol{\Omega}_0 \mathbf{X} (\mathbf{A} + \mathbf{I})]$$

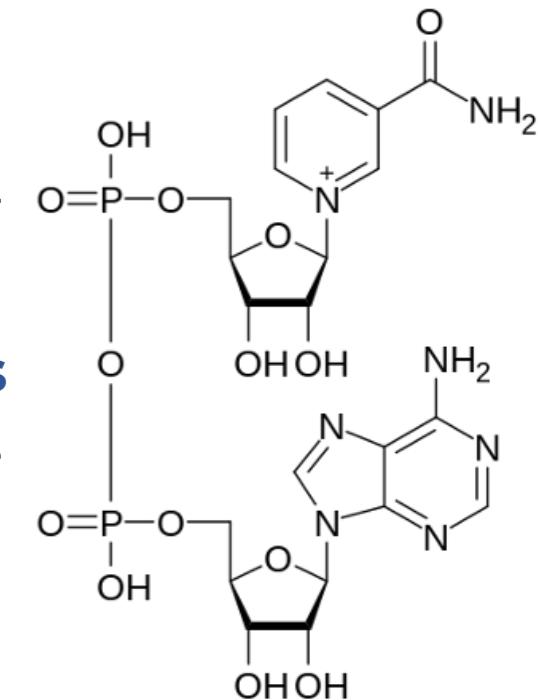
$$\mathbf{H}_2 = a[\boldsymbol{\beta}_1 \mathbf{1}^T + \boldsymbol{\Omega}_1 \mathbf{H}_1 (\mathbf{A} + \mathbf{I})]$$

$$\vdots = \vdots$$

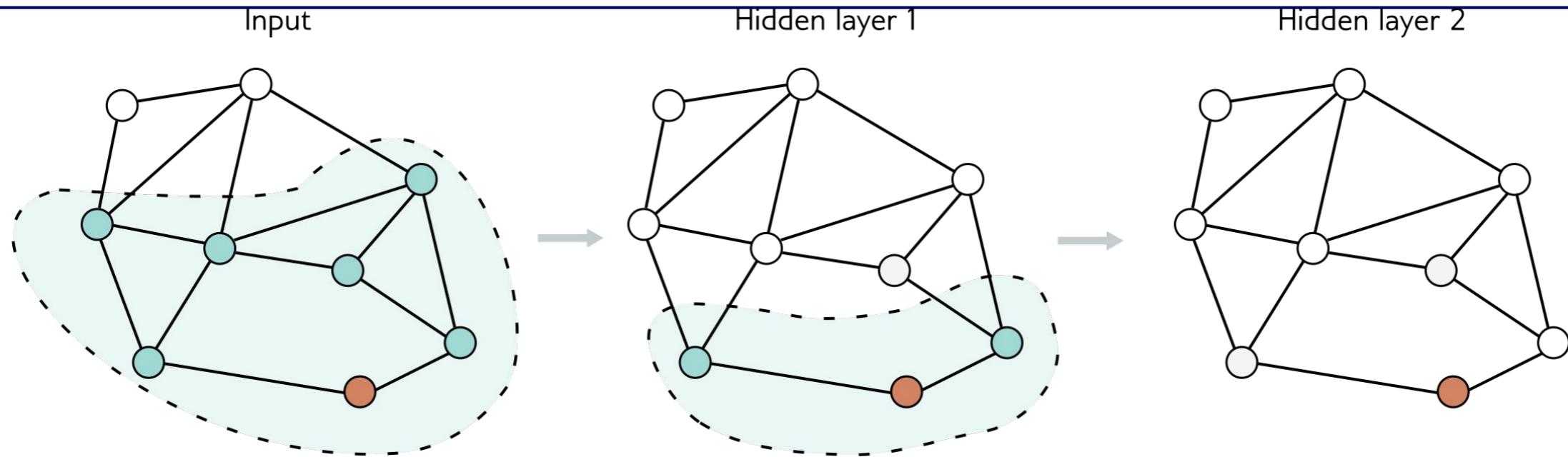
$$\mathbf{H}_k = a[\boldsymbol{\beta}_{k-1} \mathbf{1}^T + \boldsymbol{\Omega}_{k-1} \mathbf{H}_{k-1} (\mathbf{A} + \mathbf{I})]$$

$$f[\mathbf{X}, \mathbf{A}, \Phi] = \text{sig} \left[ \boldsymbol{\beta}_k + \boldsymbol{\omega}_k \mathbf{H}_k \frac{\mathbf{1}}{N} \right]$$

$$\mathbf{X}(\mathbf{A} + \mathbf{I})^k, \text{ e.g } \mathbf{H}_3: \mathbf{X}(\mathbf{A} + \mathbf{I})^3 = \mathbf{X}(\mathbf{A}^3 + 3\mathbf{A}^2 + 3\mathbf{A} + \mathbf{I}^3)$$



# GNNS (BRIEfly)



$$H^\ell \text{ depends on } X(A + I)^\ell$$

## graph expansion problem

If there are **many layers** and the graph is **densely connected**:  
every *input node* may be in the receptive field of every *output*.

In general we want that  $k \ll \text{diam}(G)$



# PART II

# EXPLAINABLE ARTIfICIAL

# INTELLIGENCE

by Bardh Prenkaj

first part of the slides based on: CSEP 590B: Explainable AI from University of Washington

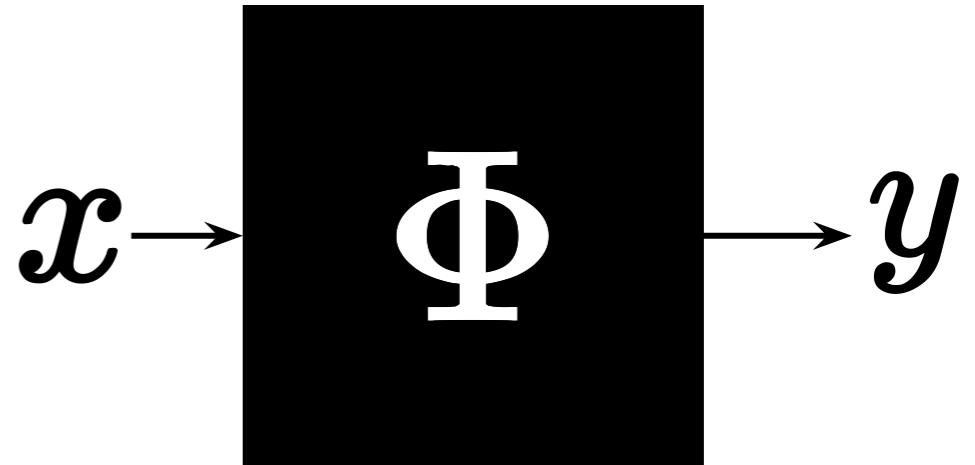


# WHAT'S GOING ON TODAY IN ML?

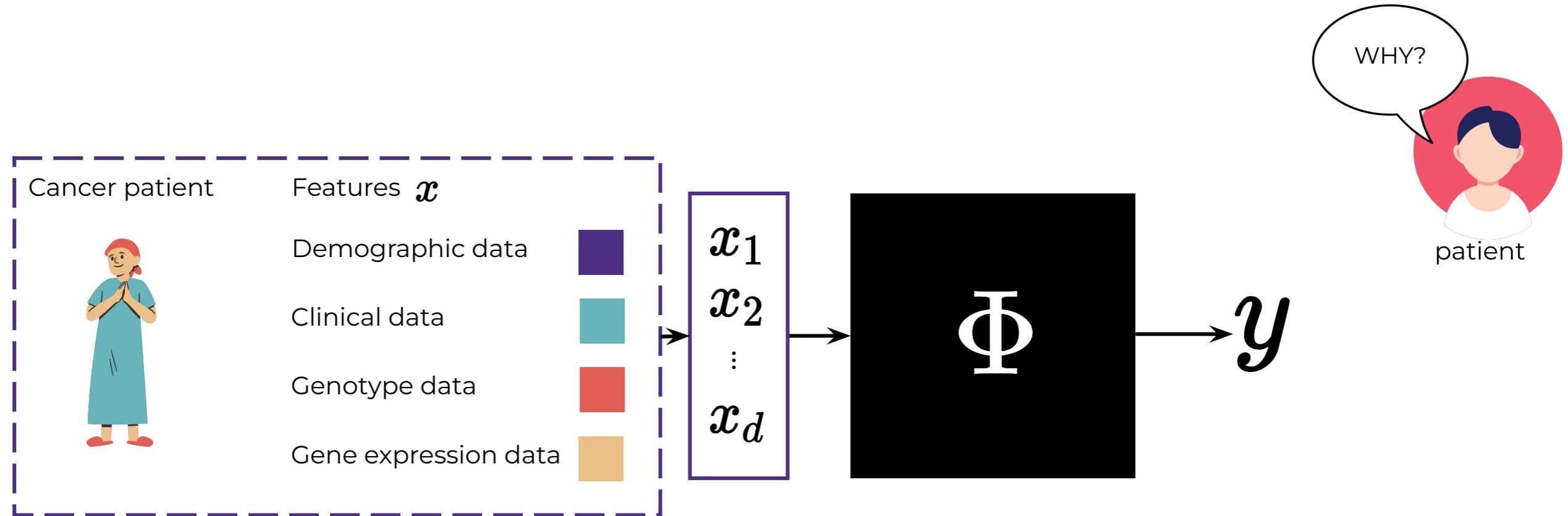
---

## Lack of transparency

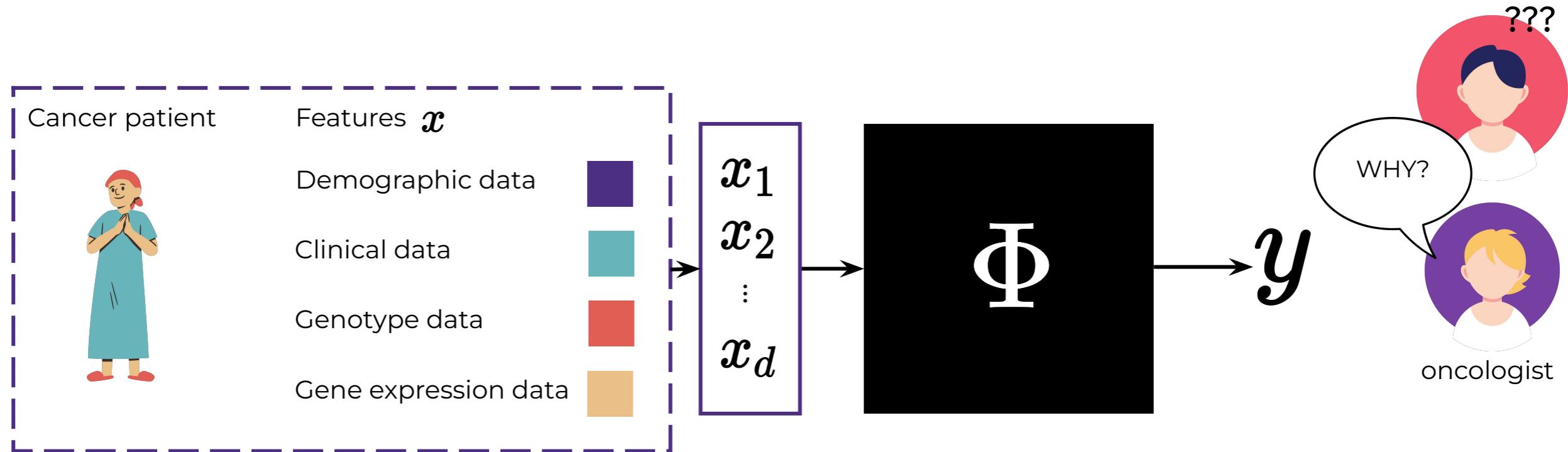
- Identify key factors in underlying processes
- Generate scientific hypotheses



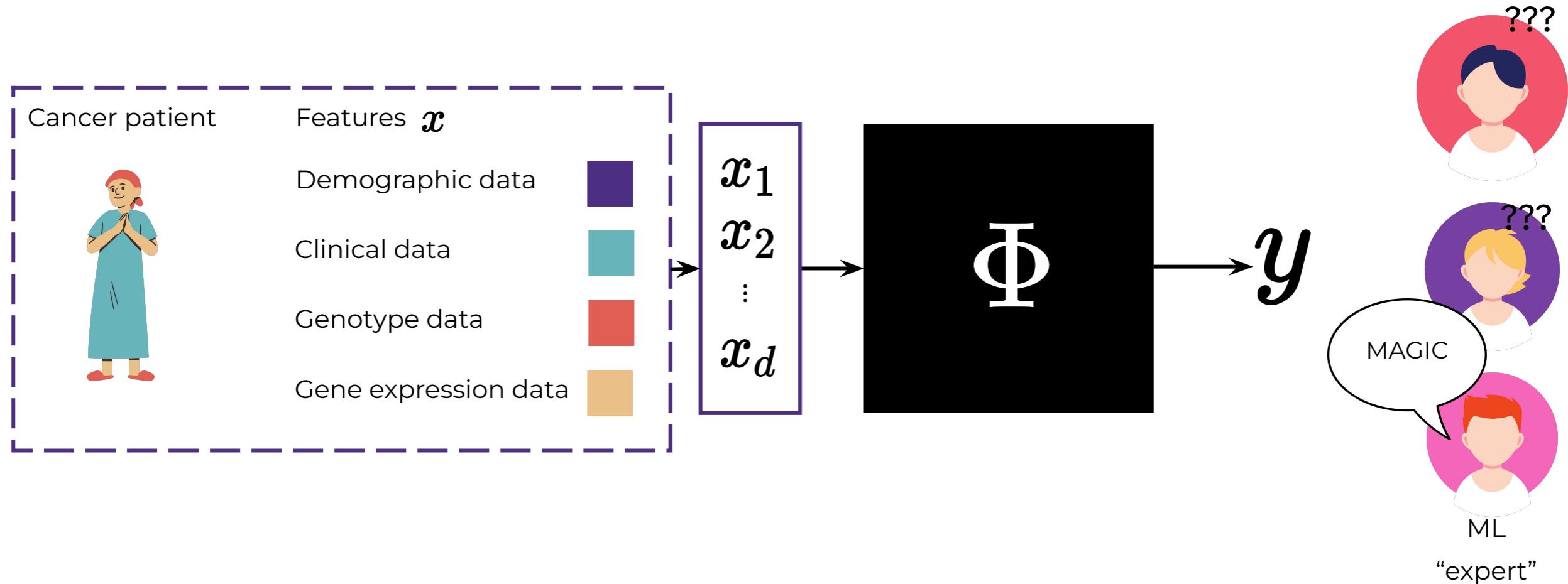
# WHY ACCURATE PREDICTIONS ARE IMPORTANT?



# WHY ACCURATE PREDICTIONS ARE IMPORTANT?



# WHY ACCURATE PREDICTIONS ARE IMPORTANT?



# EXPLAINABILITY

---

- Which features contributed to a certain prediction and how?
- How to learn or select features that are most interpretable or informative?
- How to make biological or clinical sense of a black-box model?



# TYPES OF EXPLAINABILITY

---

## Feature importance explanations

- Removal-based explanations
- Shapley values
- Propagation-based explanations

{ Some sort of factual explanations

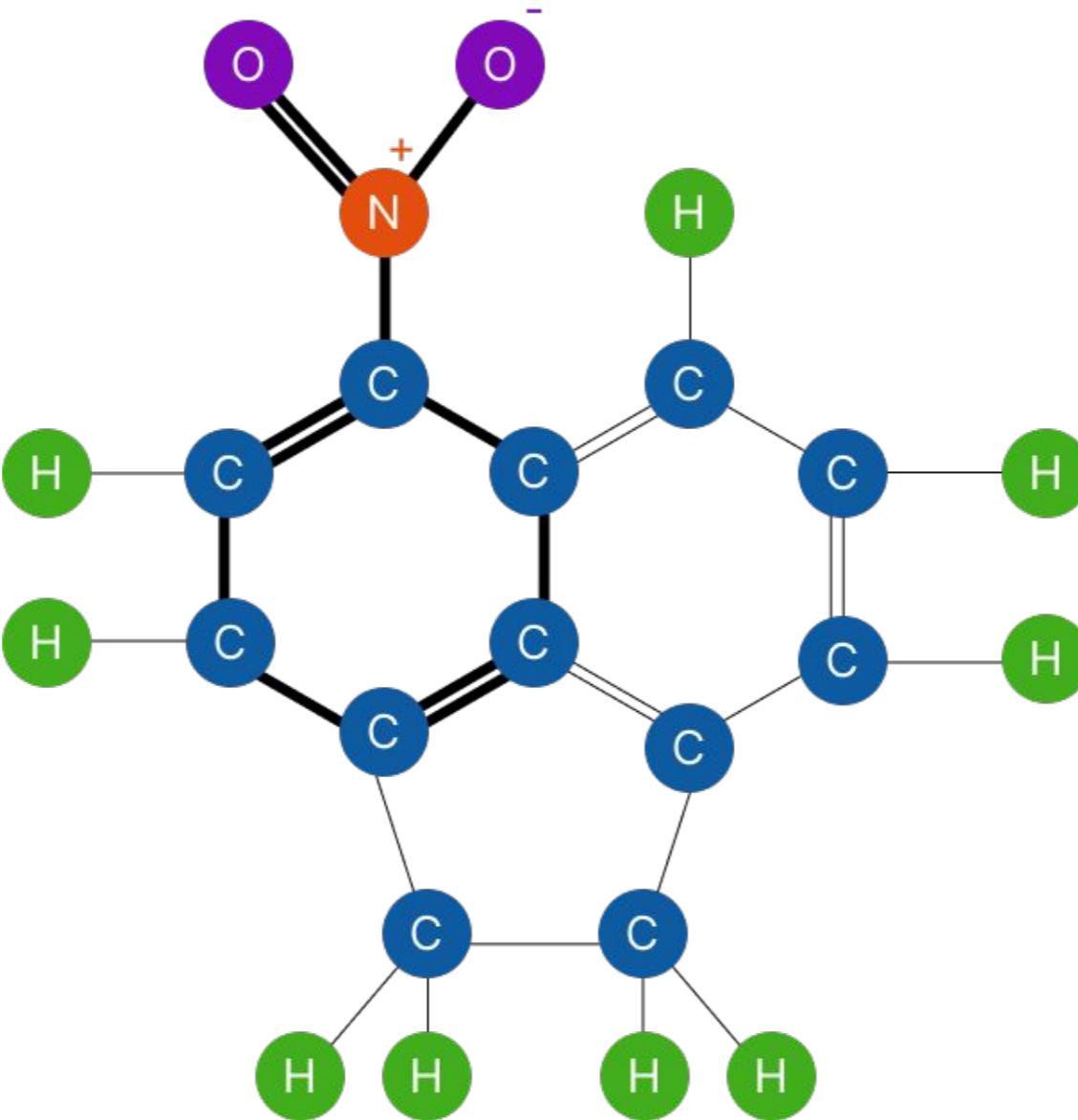
## Inherently interpretable models

## Counterfactual explanations



# FACTUAL EXPLANATIONS ON GRAPHS

---



# FACTUAL EXPLANATIONS ON GRAPHS

---

- Find those edges whose subgraph induced on them has the same label as the whole graph (**desideratum #1**)
- This subgraph should be minimal (**desideratum #2**)
- When you remove this subgraph, the remainder should have the opposite class (**corollary #1**) – *this gives sprout to factual-based counterfactual explainers*

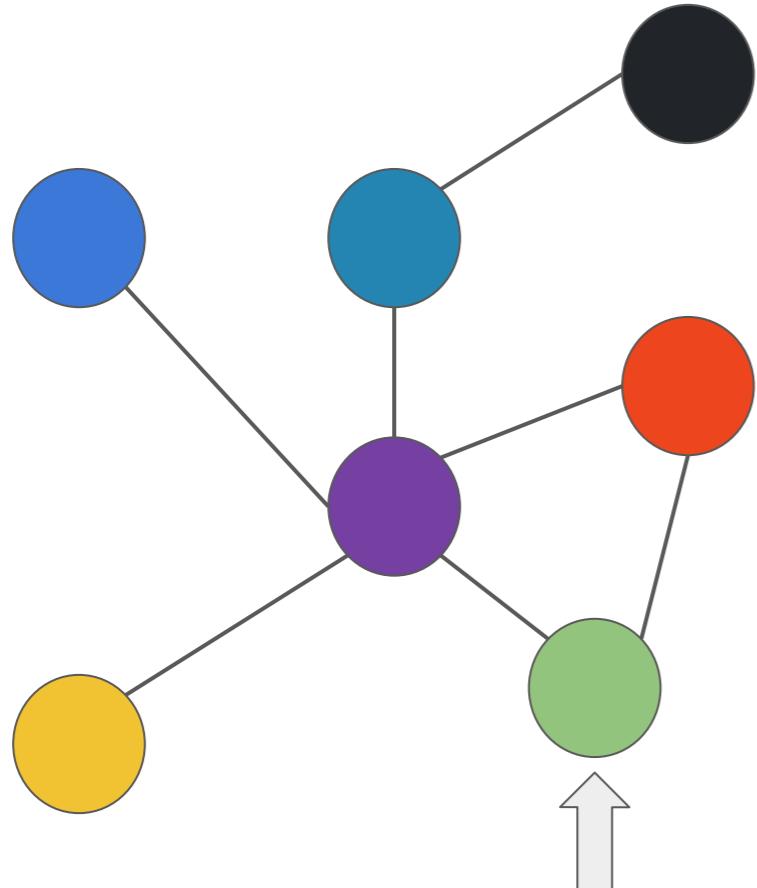




# REVISITING GNNEXPLAINER<sup>1</sup> AND GRAPHLIME<sup>2</sup>

based on (1) Ying et al. “GNNExplainer: Generating Explanations for Graph Neural Networks”, NeurIPS 2019  
& (2) Huang et al. “GraphLIME:Local Interpretable Model Explanations for Graph Neural Networks”, TKDE 2023

# GNNEPLAINER (INTRO)



node classification

## GNNEExplainer: Generating Explanations for Graph Neural Networks

Rex Ying<sup>†</sup> Dylan Bourgeois<sup>†,‡</sup> Jiaxuan You<sup>†</sup> Marinka Zitnik<sup>†</sup> Jure Leskovec<sup>†</sup>

<sup>†</sup>Department of Computer Science, Stanford University

<sup>‡</sup>Robust.AI

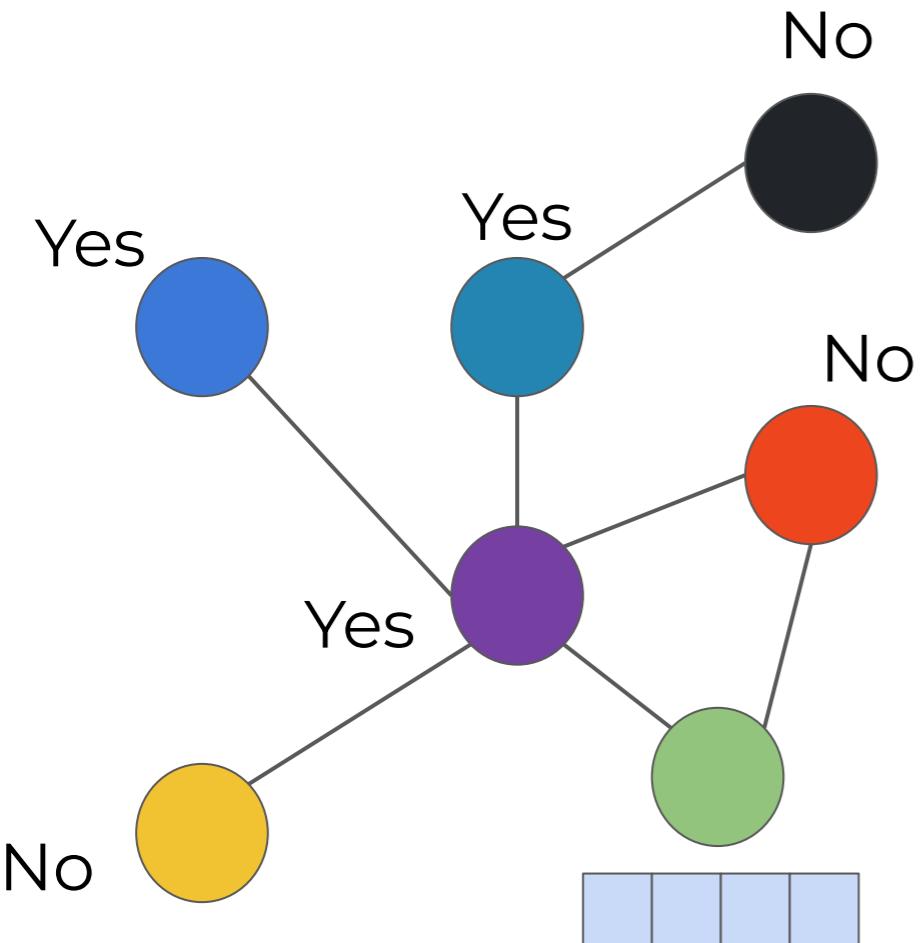
{rexying, dtsbourg, jiaxuan, marinka, jure}@cs.stanford.edu

### Abstract

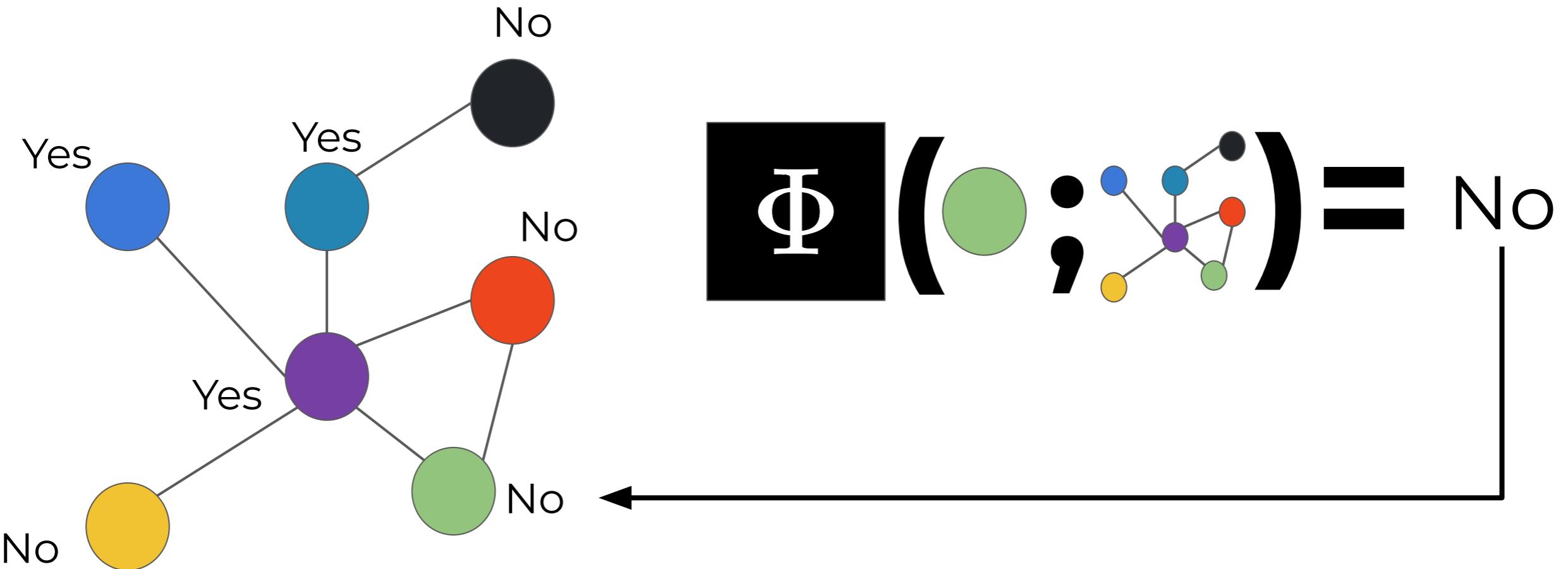
Graph Neural Networks (GNNs) are a powerful tool for machine learning on graphs. GNNs combine node feature information with the graph structure by recursively passing neural messages along edges of the input graph. However, incorporating both graph structure and feature information leads to complex models and explaining predictions made by GNNs remains unsolved. Here we propose GNNEPLAINER, the first general, model-agnostic approach for providing interpretable explanations for predictions of any GNN-based model on any graph-based machine learning task. Given an instance, GNNEPLAINER identifies a compact subgraph structure and a small subset of node features that have a crucial role in GNN's prediction. Further, GNNEPLAINER can generate consistent and concise explanations for an entire class of instances. We formulate GNNEPLAINER as an optimization task that maximizes the mutual information between a GNN's prediction and distribution of possible subgraph structures. Experiments on synthetic and real-world graphs show that our approach can identify important graph structures as well as node features, and outperforms alternative baseline approaches by up to 43.0% in explanation accuracy. GNNEPLAINER provides a variety of benefits, from the ability to visualize semantically relevant structures to interpretability, to giving insights into errors of faulty GNNs.

# GNNEXPLAINER (INTRO)

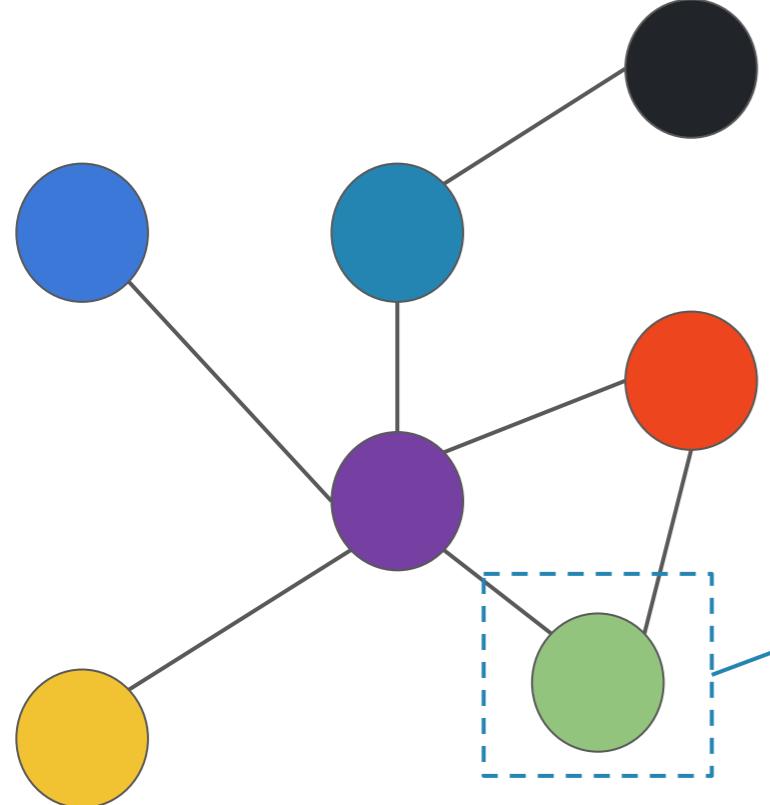
---



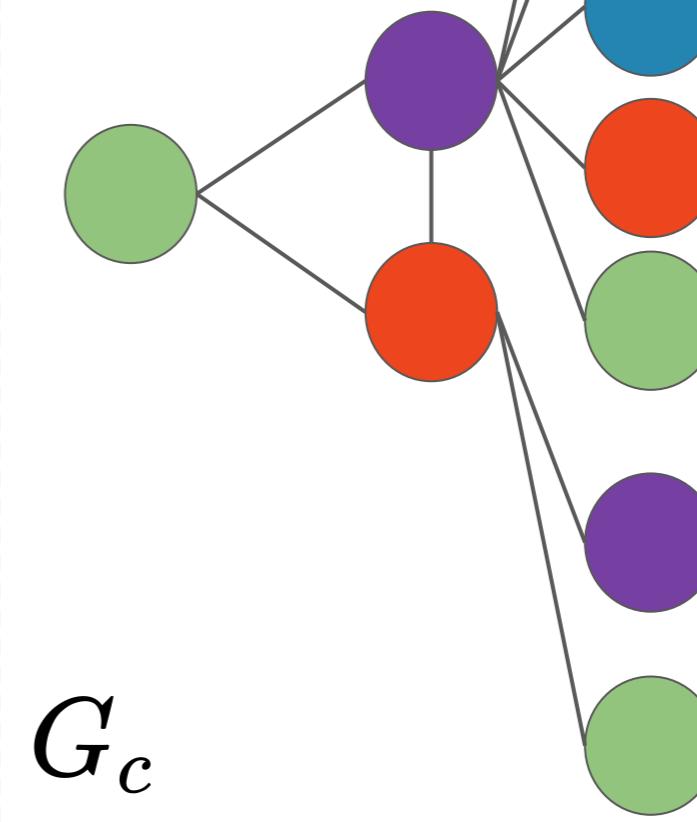
# GNNEXPLAINER (INTRO)



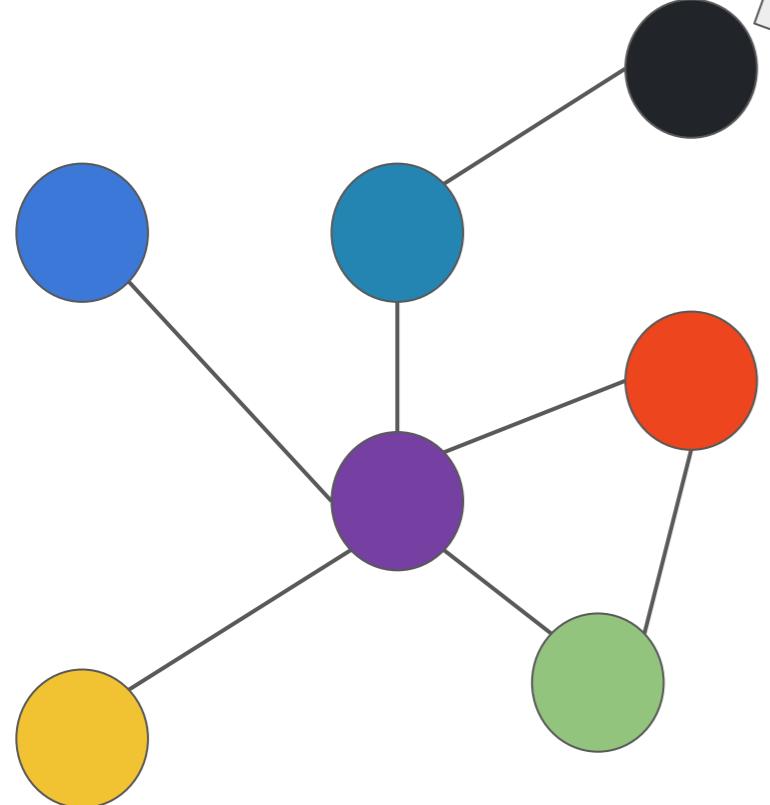
# GNNEXPLAINER (COMP. GRAPH)



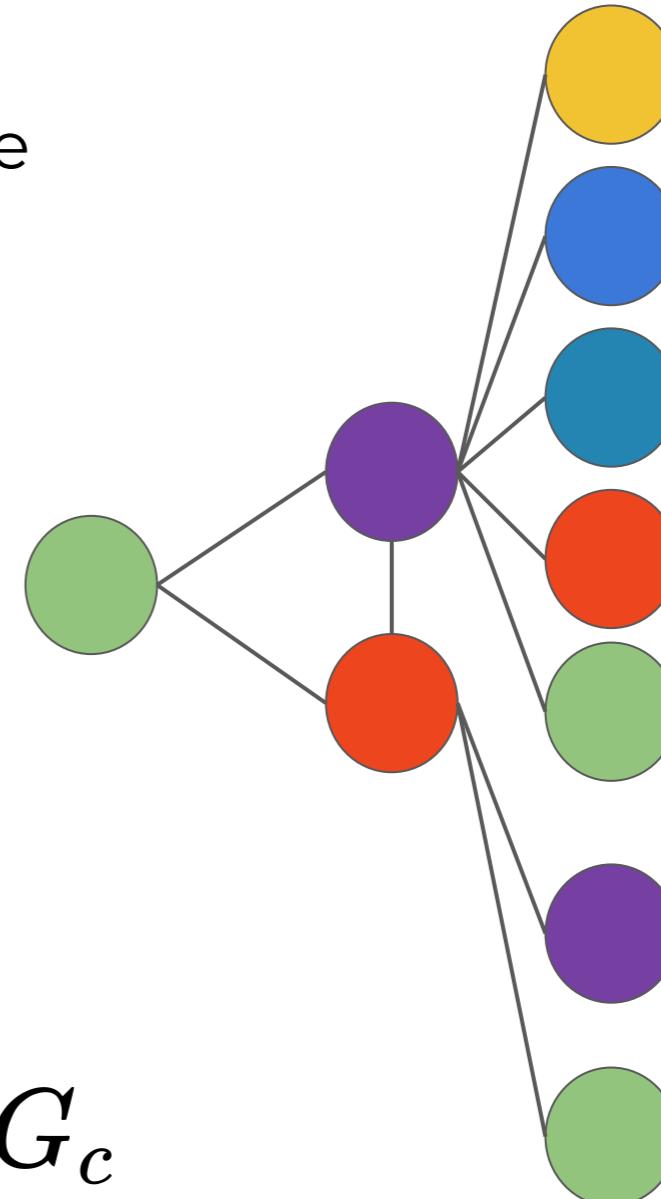
computational  
graph for a  
2-layer GNN



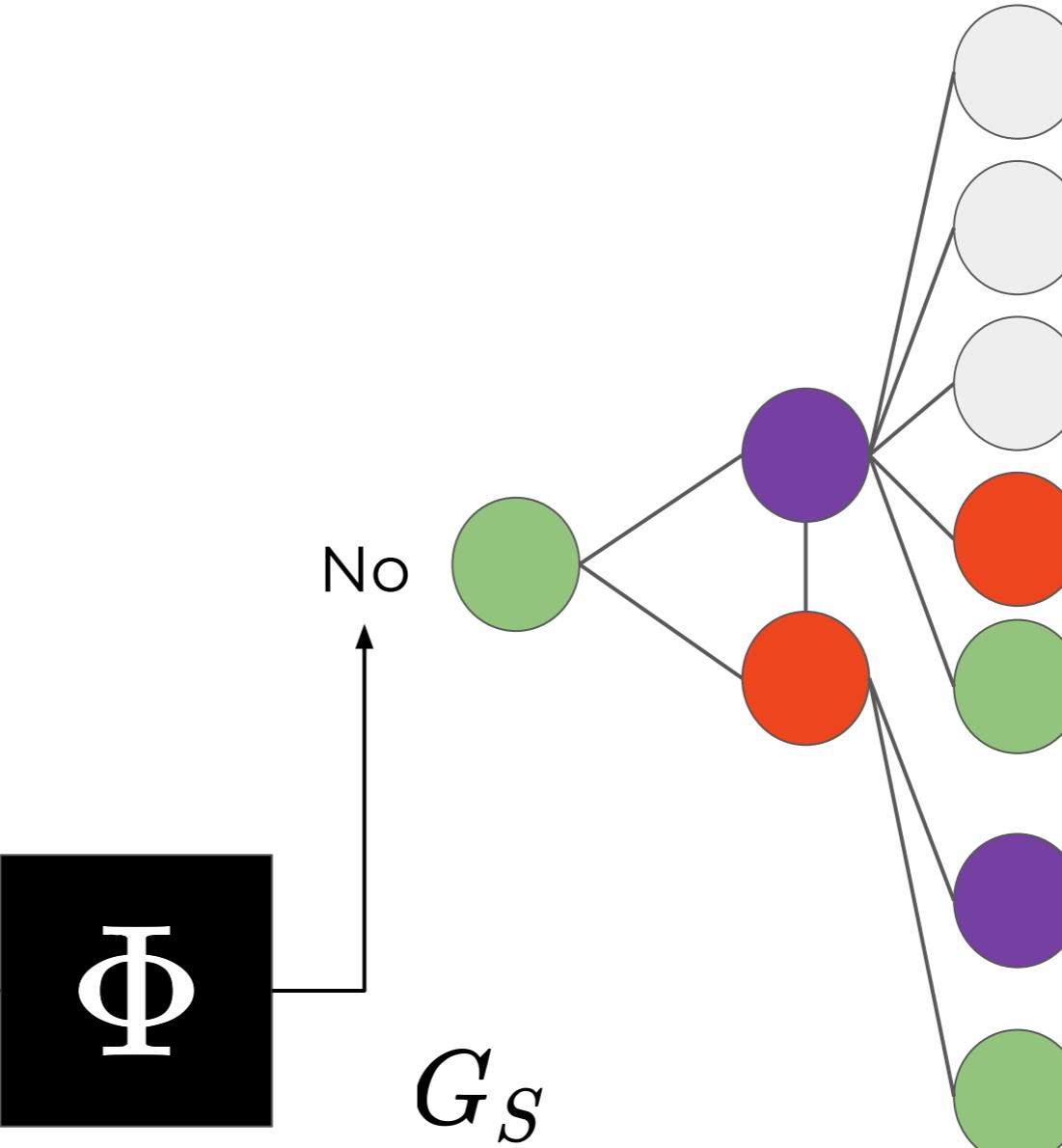
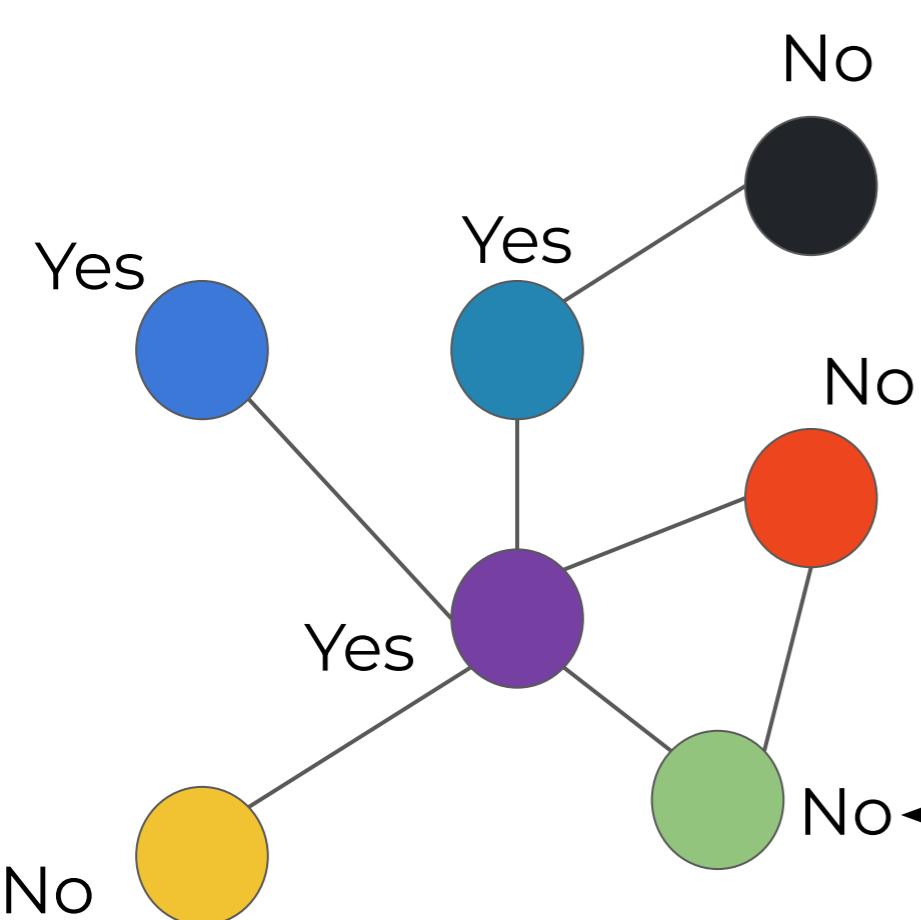
# GNNEXPLAINER (COMP. GRAPH)



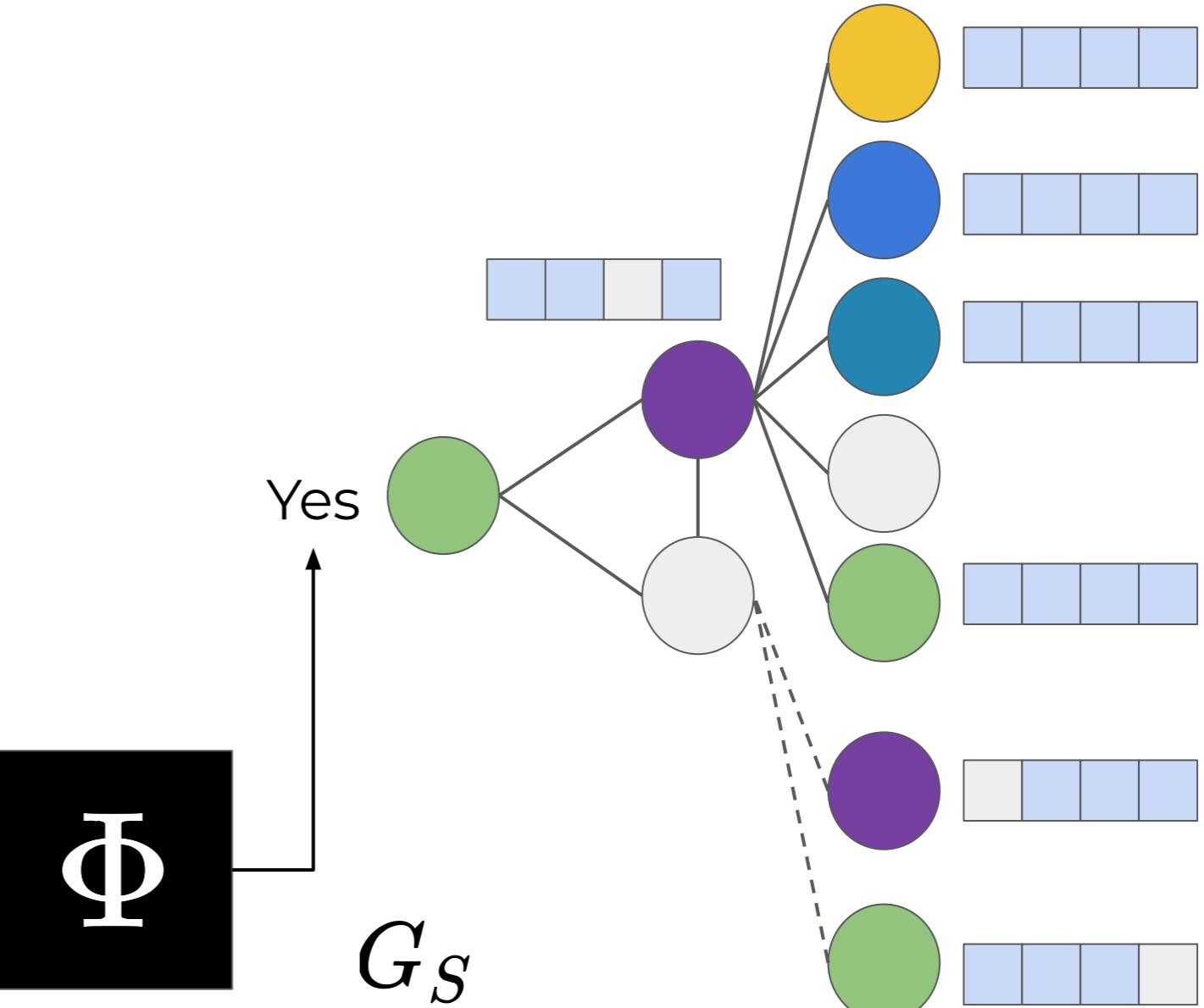
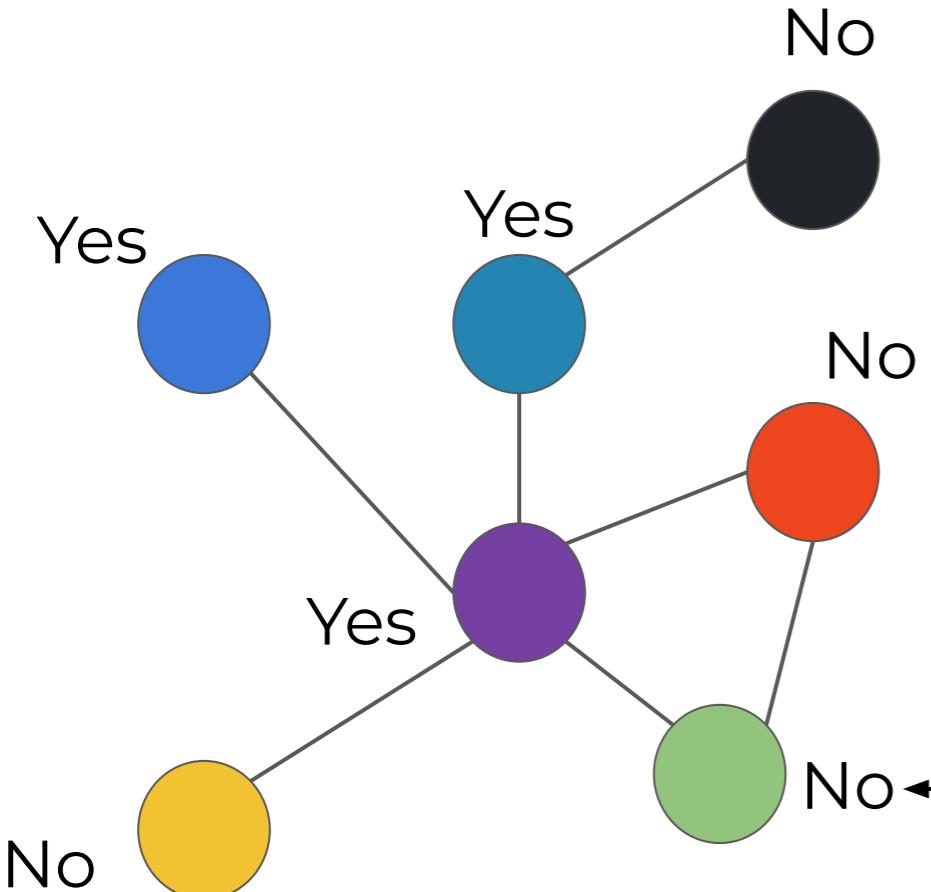
out of scope since  
it's too far away



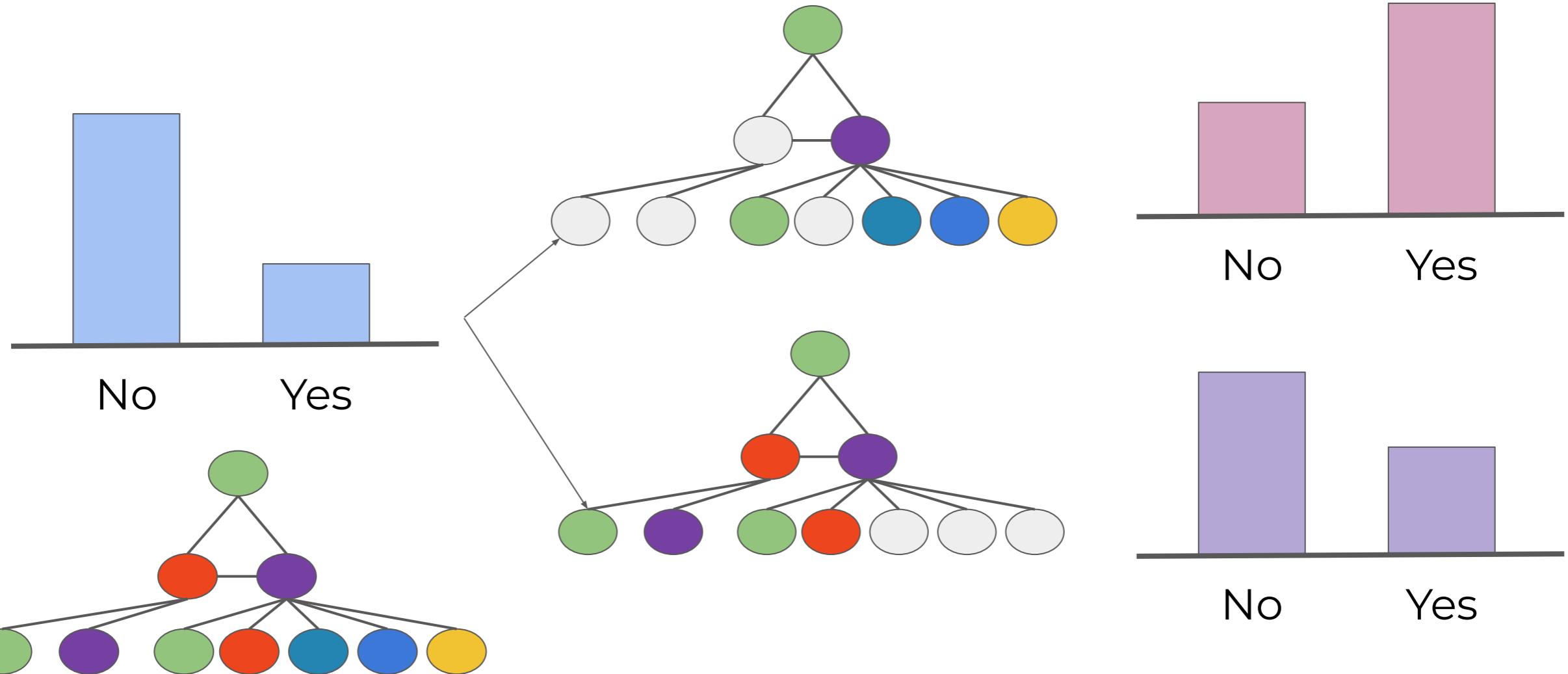
# GNNEXPLAINER (REMOVE NODES)



# GNNEXPLAINER (REMOVE NODES)



# GNNEXPLAINER (PROBABILITY DISTRIBUTIONS)



# GNNEXPLAINER (FORMAL)

---

$$\max_{G_S} MI(Y, (G_S, X_S)) = H(Y) - H(Y \mid G = G_S, X = X_S)$$

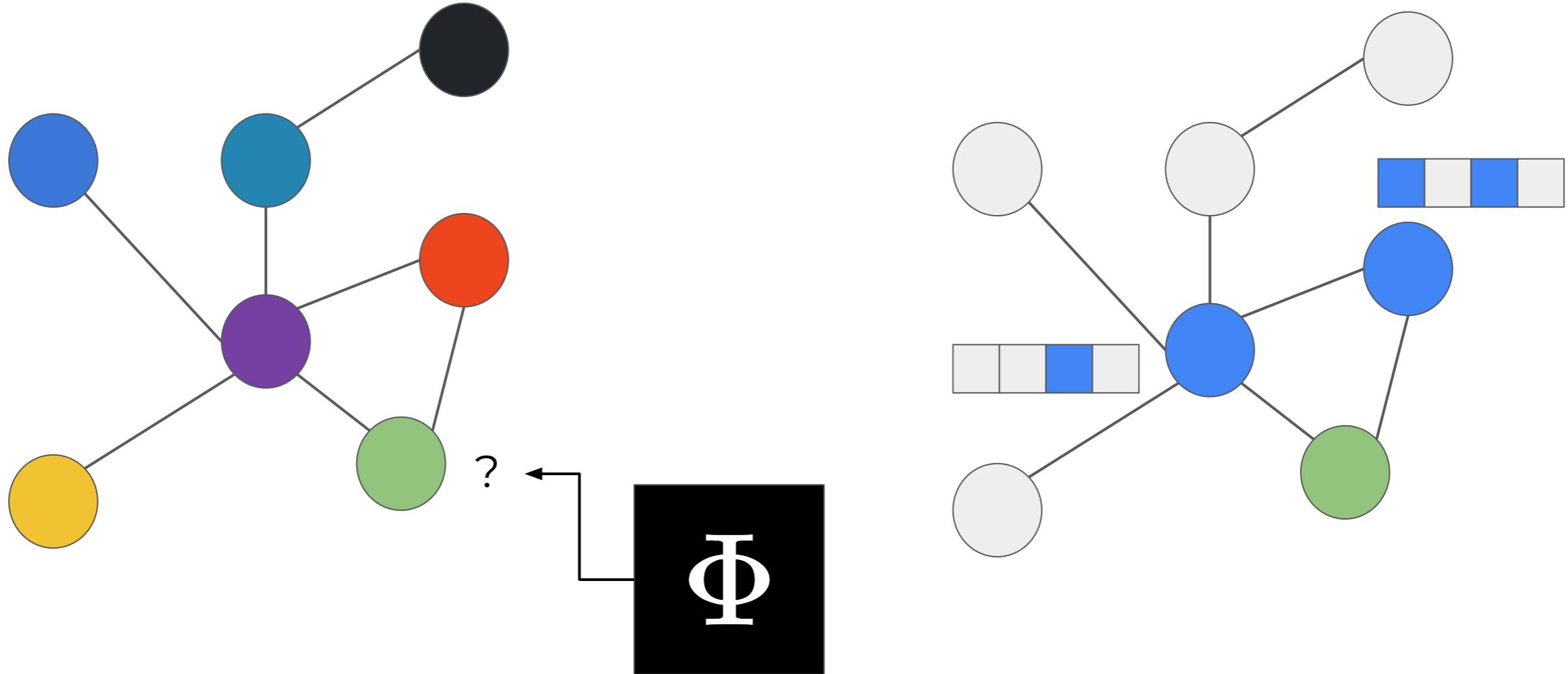
## Human-interpretable intuition

Find a subgraph that maintains as much information as possible compared to the full graph



# GNNEXPLAINER (EXPLANATION EXAMPLE)

---



# GRAPHLIME IN A NUTSHELL

---

- GraphLIME builds **local interpretable models** for specific nodes in the graph
- It uses the **Hilbert-Schmidt Independence Criterion (HSIC) Lasso** for feature selection.
- GraphLIME operates **locally within the subgraph** of the node being explained.

M. Yamada et al., “Ultra high-dimensional nonlinear feature selection for big biological data,” IEEE Trans. Knowl. Data Eng., vol. 30, no. 7, pp. 1352–1365, Jul. 2018.



# SELECTING THE NODE'S NEIGHBORHOOD

---

- For a target node, GraphLIME selects its **N-hop neighborhood** (where N is the number of GNN layers)
- This neighborhood captures relevant **structural information** around the node



# BUILDING THE LOCAL MODEL

---

- GraphLIME learns a **nonlinear interpretable model** within the selected subgraph
- This model explains the node's **prediction behavior**
- GraphLIME computes the **K most representative features** using HSIC Lasso



# FEATURE IMPORTANCE WITH HSIC LASSO

---

- HSIC Lasso is a **nonlinear method** for selecting important features
- It balances the trade-off between **model complexity and interpretability**
- HSIC Lasso identifies the most influential features for the node's prediction





# PART III

# COUNTERFACTUAL

# EXPLANATIONS IN GRAPHS

by Mario A. Prado-Romero

Based on:

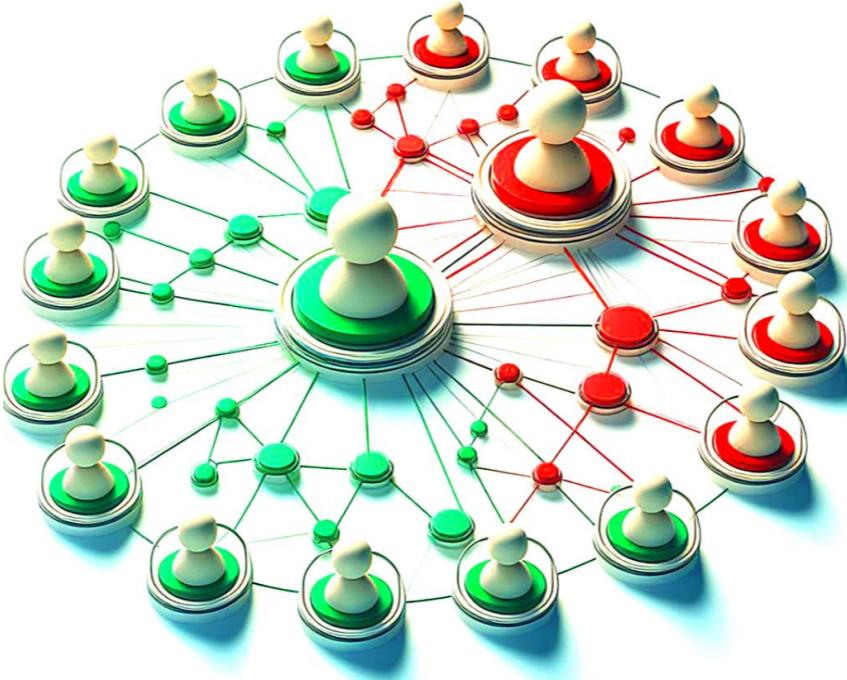
Prado-Romero et al. "[A survey on graph counterfactual explanations: definitions, methods, evaluation](#)", ACM CSUR 2023



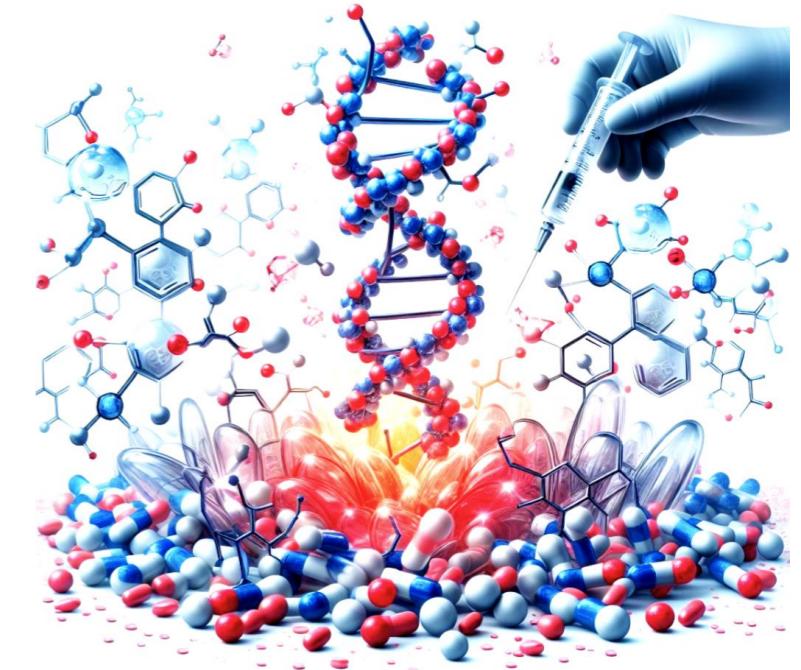
# GRAPH COUNTERFACTUALS?



Road design in  
Transportation Networks

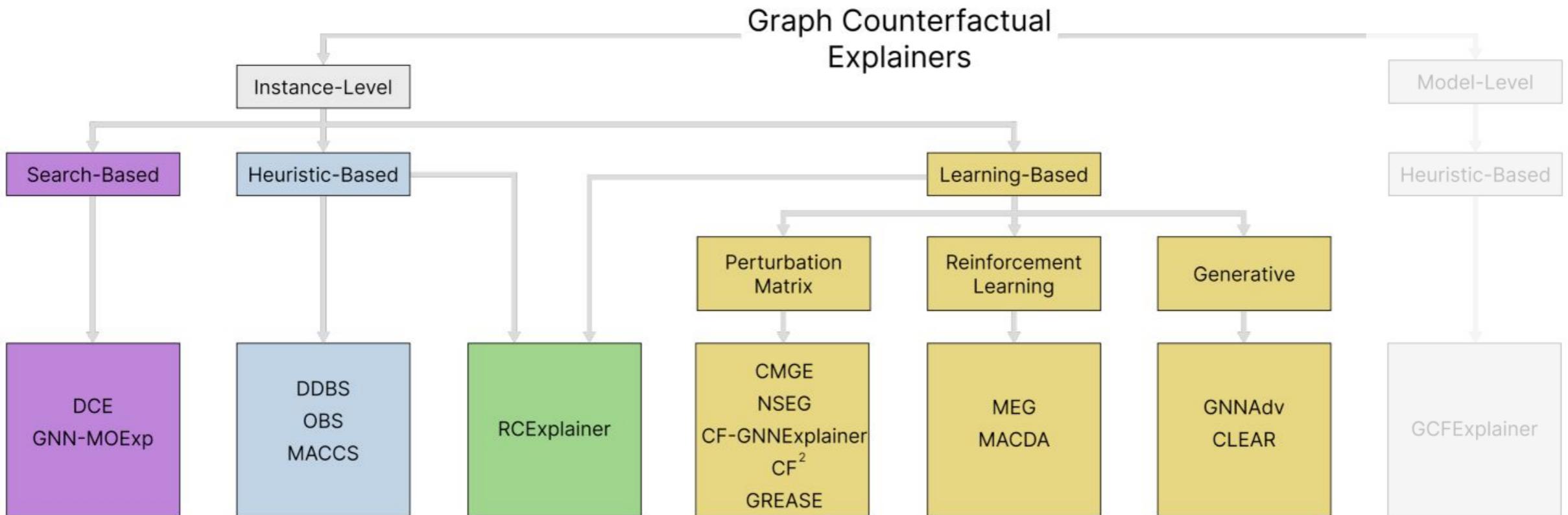


Recommend potential  
collaborations in social  
networks to reach a target  
audience



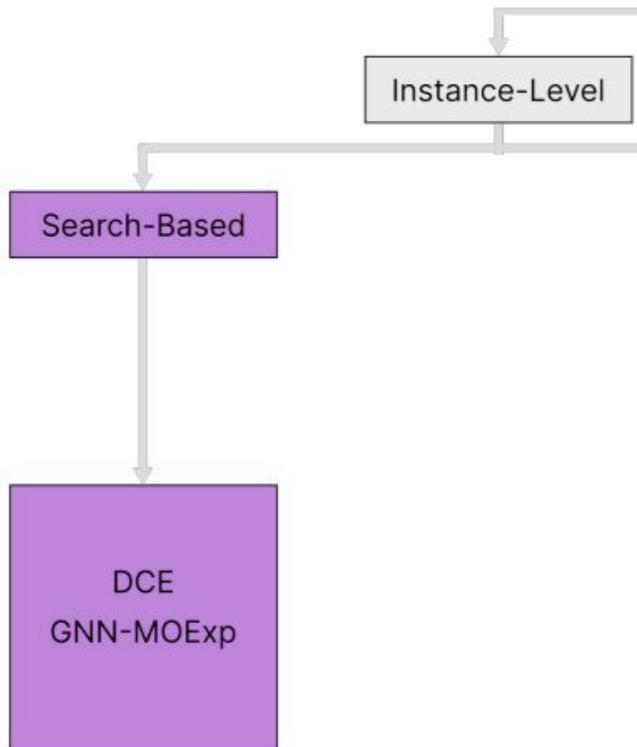
Modifying molecular  
structures for drug discovery

# GCE METHODS



# SEARCH-BASED GCE METHODS

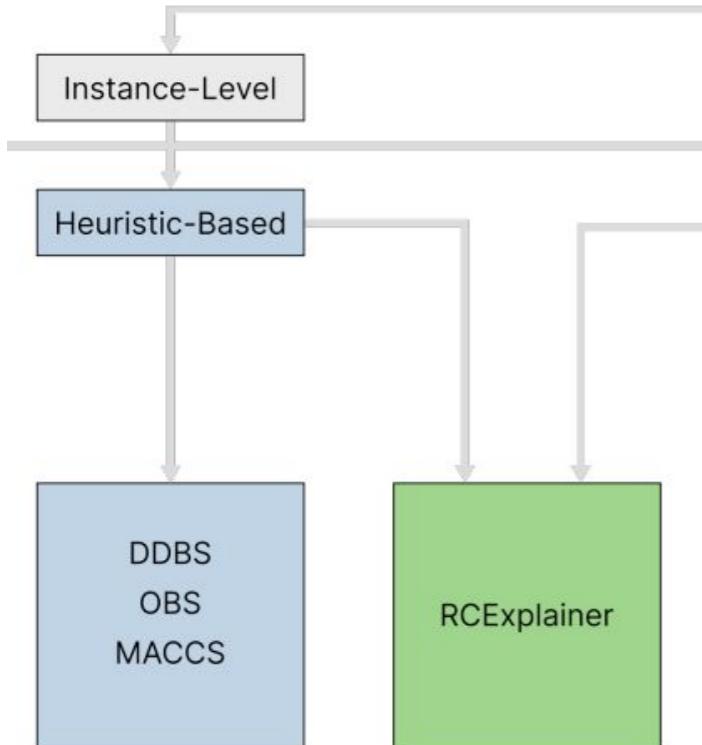
---



- Find a counterfactual **within the data**
- For a graph  $\mathbf{G} \in \mathcal{G}$  find a  $\mathbf{G}' \in \mathcal{G}$  s.t.  $\Phi(\mathbf{G}) \neq \Phi(\mathbf{G}')$
- These methods **fail** to produce a counterfactual if the explainer **cannot access** the original **dataset**

# HEURISTIC-BASED GCE METHODS

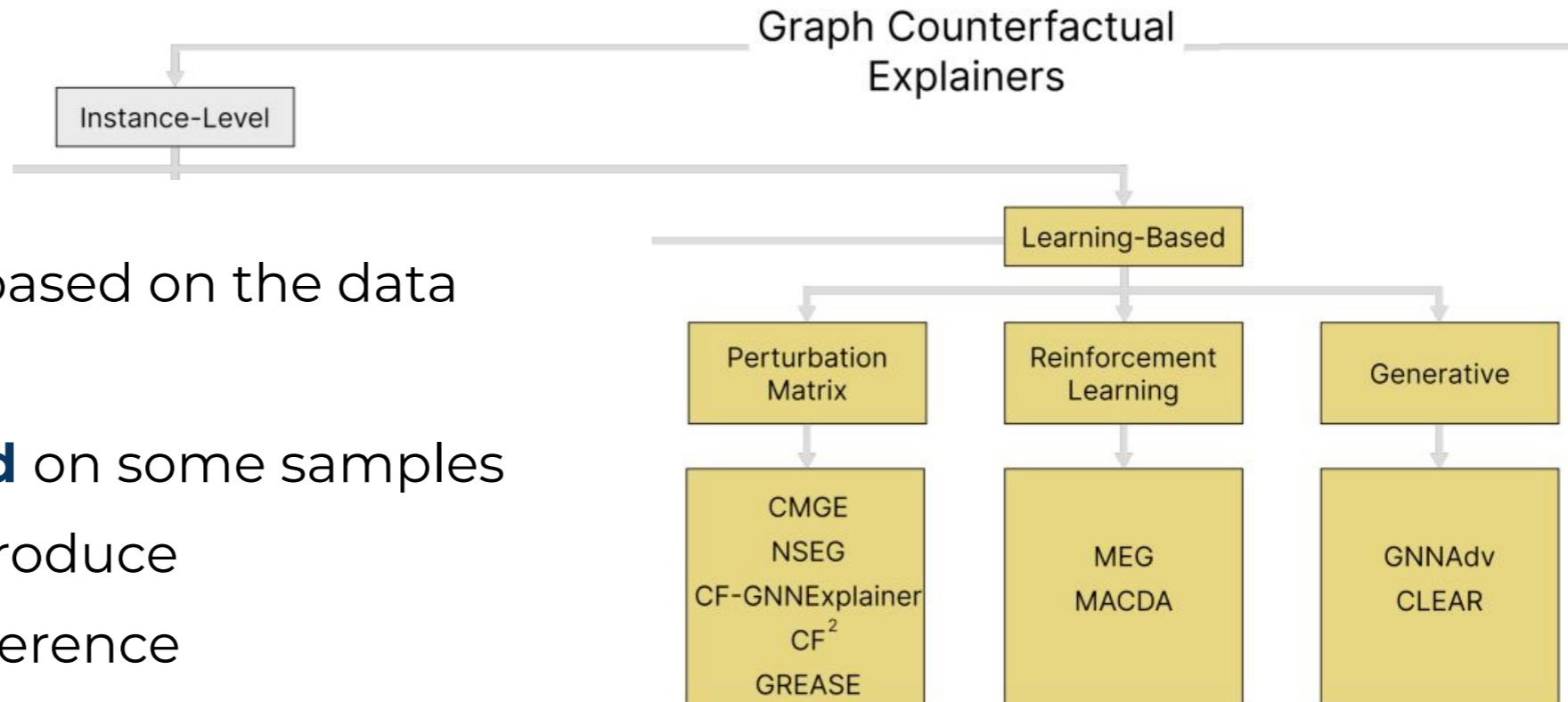
---



- **Perturb** the original graph such that  $\Phi(\mathbf{G}) \neq \Phi(\mathbf{G}')$  **without** accessing the original dataset
- **Requires** to define the perturbation **rules** after a **careful examination** of the data

# LEARNING-BASED GCE METHODS

- **Learn the heuristic** based on the data
- Explainers are **trained** on some samples and can be used to produce counterfactuals at inference

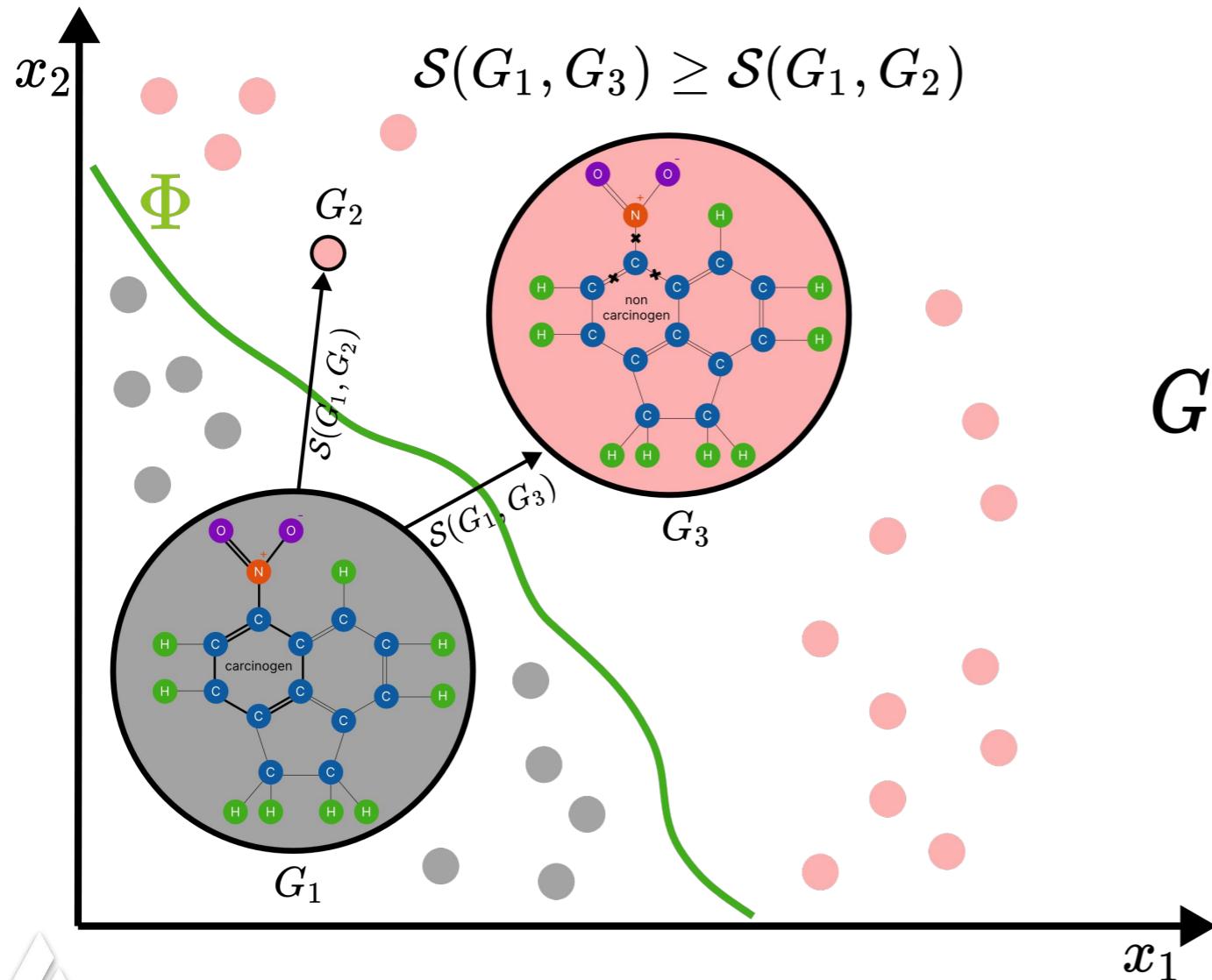




# BASELINE GRAPH COUNTERFACTUAL EXPLAINER

L. Faber, A. K. Moghaddam, and R. Wattenhofer. 2020. Contrastive Graph Neural Network Explanation. In Proc. of the 37th Graph Repr. Learning and Beyond Workshop at ICML 2020. Int. Conf. on Machine Learning, 28

# DCE (DISTRIBUTION COMPLAINT EXPLANATIONS)



$$G^* = \arg \min_{G' \in \mathcal{G}, \Phi(G) \neq \Phi(G')} d(G, G')$$

# QUESTIONS ON DCE

---

- Does it guarantee to always produce valid counterfactuals?
- Does it guarantee to produce the closest counterfactual to the input we want to explain?
- Are counterfactuals within the data manifold?





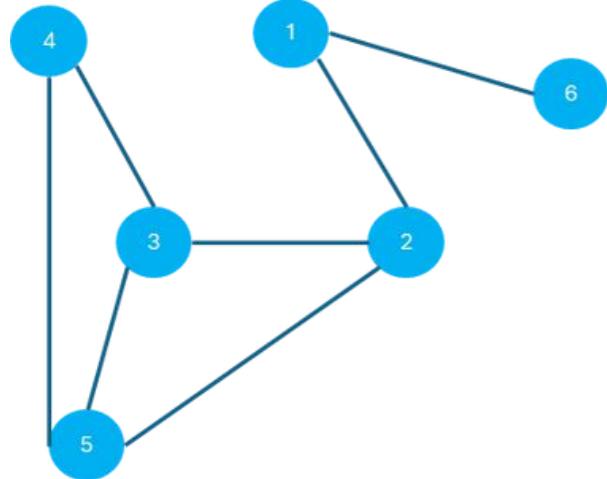
# HEURISTIC-BASED EXPLAINERS

Abbate C, Bonchi F. Counterfactual graphs for explainable classification of brain networks. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining 2021 Aug 14 (pp. 2495-2504).

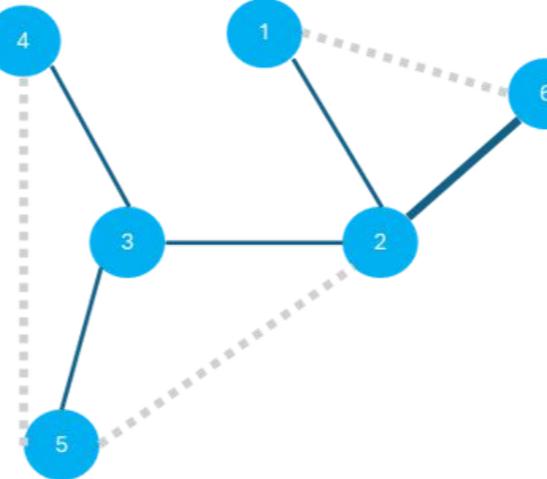
Wellawatte GP, Gandhi HA, Seshadri A, White AD. A Perspective on Explanations of Molecular Prediction Models. Journal of Chemical Theory and Computation. 2023 Mar 27;19(8):2149-60.

# OBLIVIOUS BIDIRECTIONAL SEARCH (OBS)

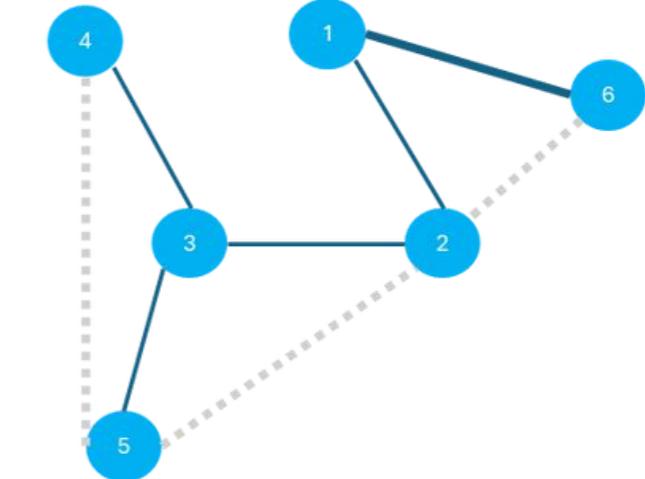
---



Cyclic  
Graph

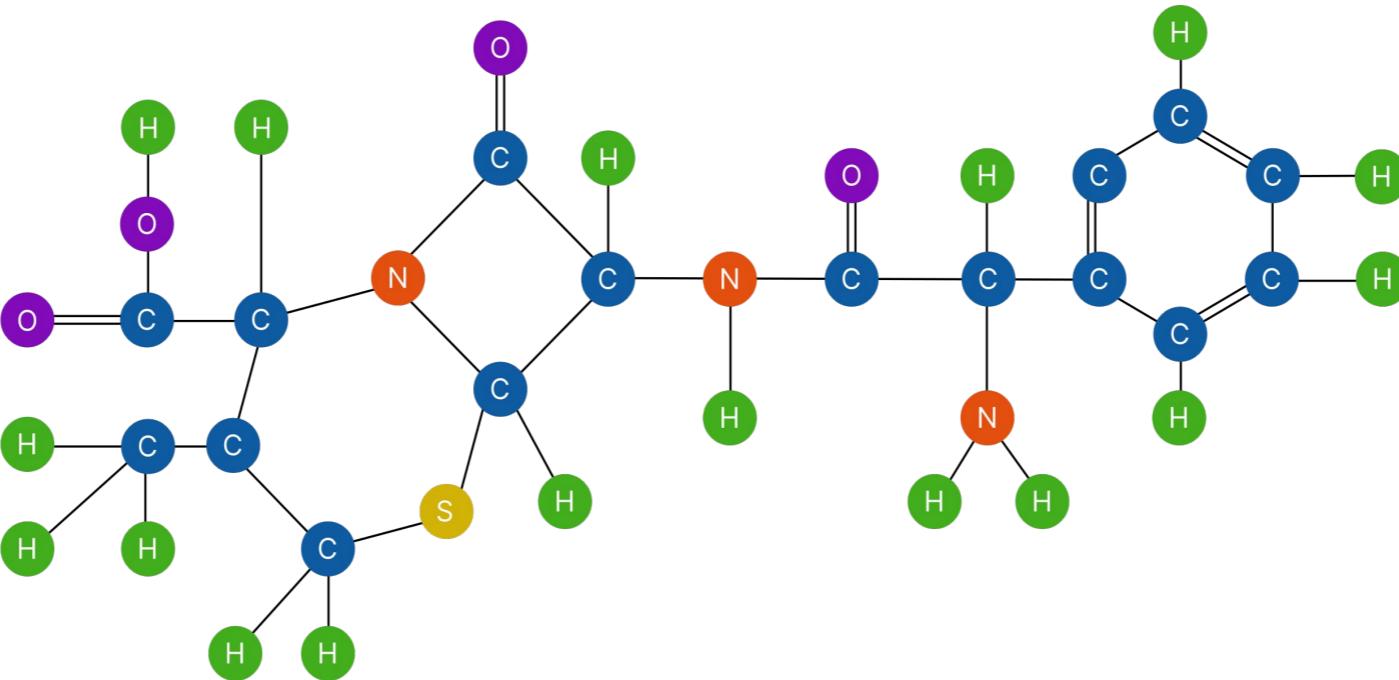


Step 1:  
Find a  
Counterfactual



Step 2:  
Reduce distance between  
original graph and  
counterfactual

Uses SMILES representations of molecules



**CC1=C(N2C(C(C2=O)NC(=O)C(C3=CC=CC=C3)N)SC1)C(=O)O**



# GLOBAL COUNTERFACTUAL EXPLAINERS (DIGRESSION)

Zexi Huang, Mert Kosan, Sourav Medya, Sayan Ranu, and Ambuj Singh. 2023. Global Counterfactual Explainer for Graph Neural Networks. In Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining (WSDM '23). Association for Computing Machinery, New York, NY, USA, 141–149. <https://doi.org/10.1145/3539597.3570376>

# GCFEXPLAINER

---

- The **input** is a set of instances to explain  $\mathbb{G}$
- The **output** is a set of counterfactual instances  $\mathbb{C}$  that are considered an explanation of the input set
- **Coverage:** the proportion of graphs in  $\mathbb{G}$  that have a close counterfactual in  $\mathbb{C}$  under a given distance threshold  $\theta$ :

$$\text{coverage}(\mathbb{C}) = |\{G \in \mathbb{G} \mid \min_{\{C \in \mathbb{C}\}} \{d(G, C)\} \leq \theta\}| / |\mathbb{G}|$$



# SEARCH SPACE

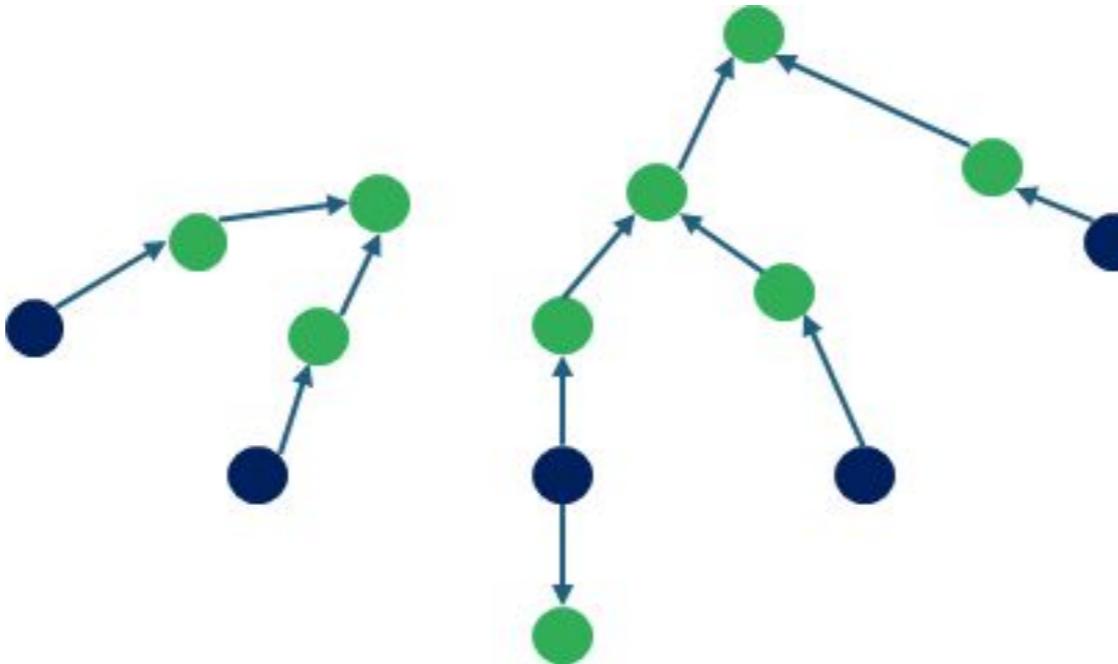
---

- The search space of counterfactual graphs = **graphs in the same domain as the input within a distance of  $\theta$**
- The number of potential graphs within  $\theta$  **increases exponentially** since the space of graph edits is combinatorial
- Use an edit map to organize these graphs as a meta-graph ⚡



# SEARCH SPACE (ILLUSTRATIVE)

---



# VERTEX-REINFORCED RANDOM WALK

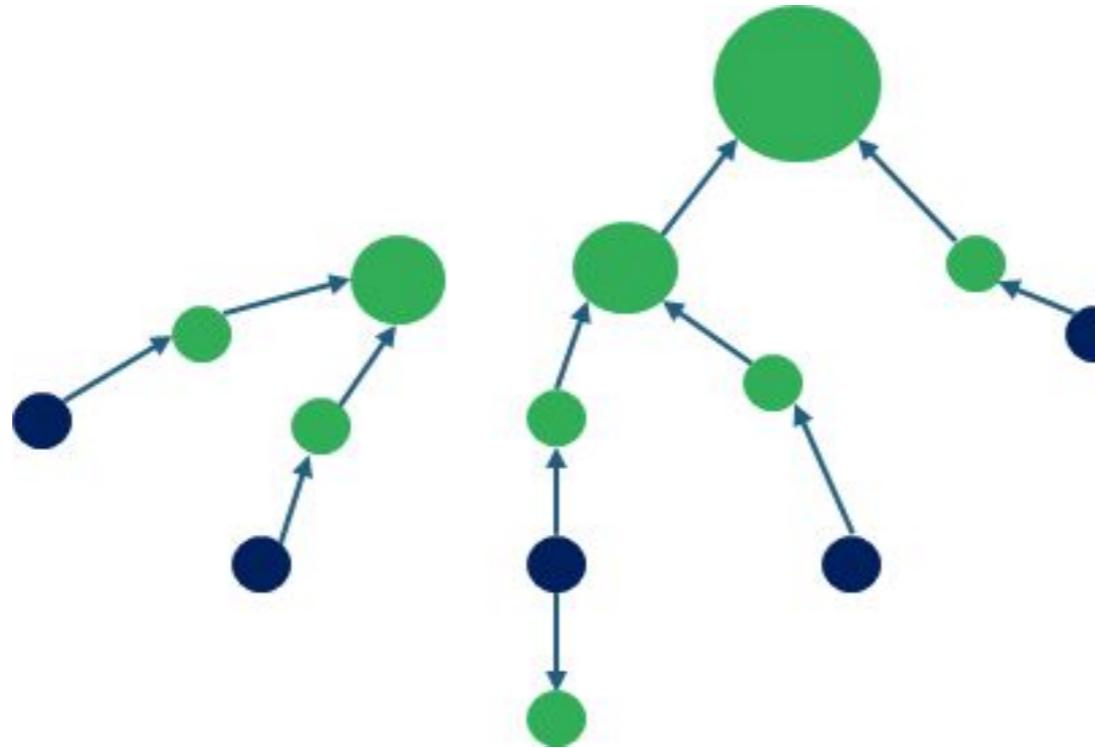
---

- Leverage vertex-reinforced random walks (VRRW) on the edit map  $\mathfrak{G}$
- VRRW converges to a set of nodes that are **important** and **diverse**, which will form a small set of counterfactual candidates for further processing



# VERTEX-REINFORCED RANDOM WALK

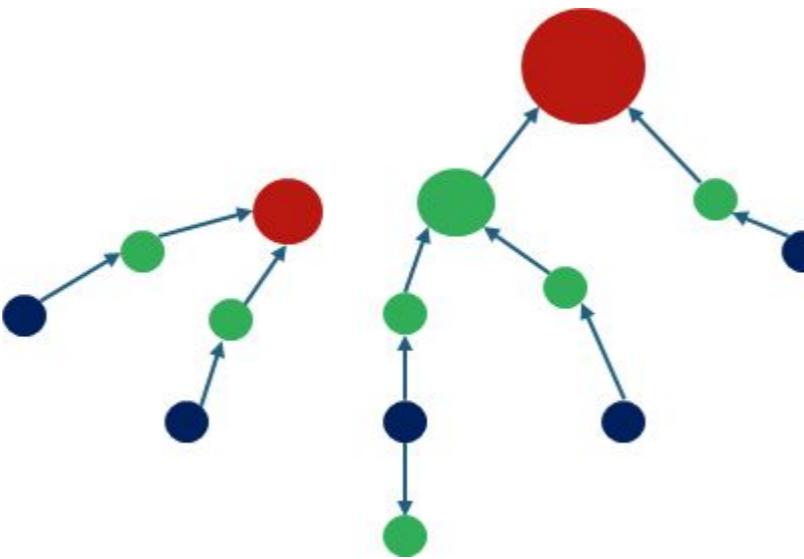
---



# ITERATIVE SUMMARY COMPUTATION

---

Create a summary of counterfactual graphs as the recourse representation by iteratively adding the best candidate based on the **maximal gain of the coverage**





# LEARNING-BASED EXPLAINERS

D. Numeroso and D. Bacci. 2021. Meg: Generating molecular counterfactual explanations for deep graph networks. In 2021 Int. Joint Conf. on Neural Networks. IEEE, 1–8

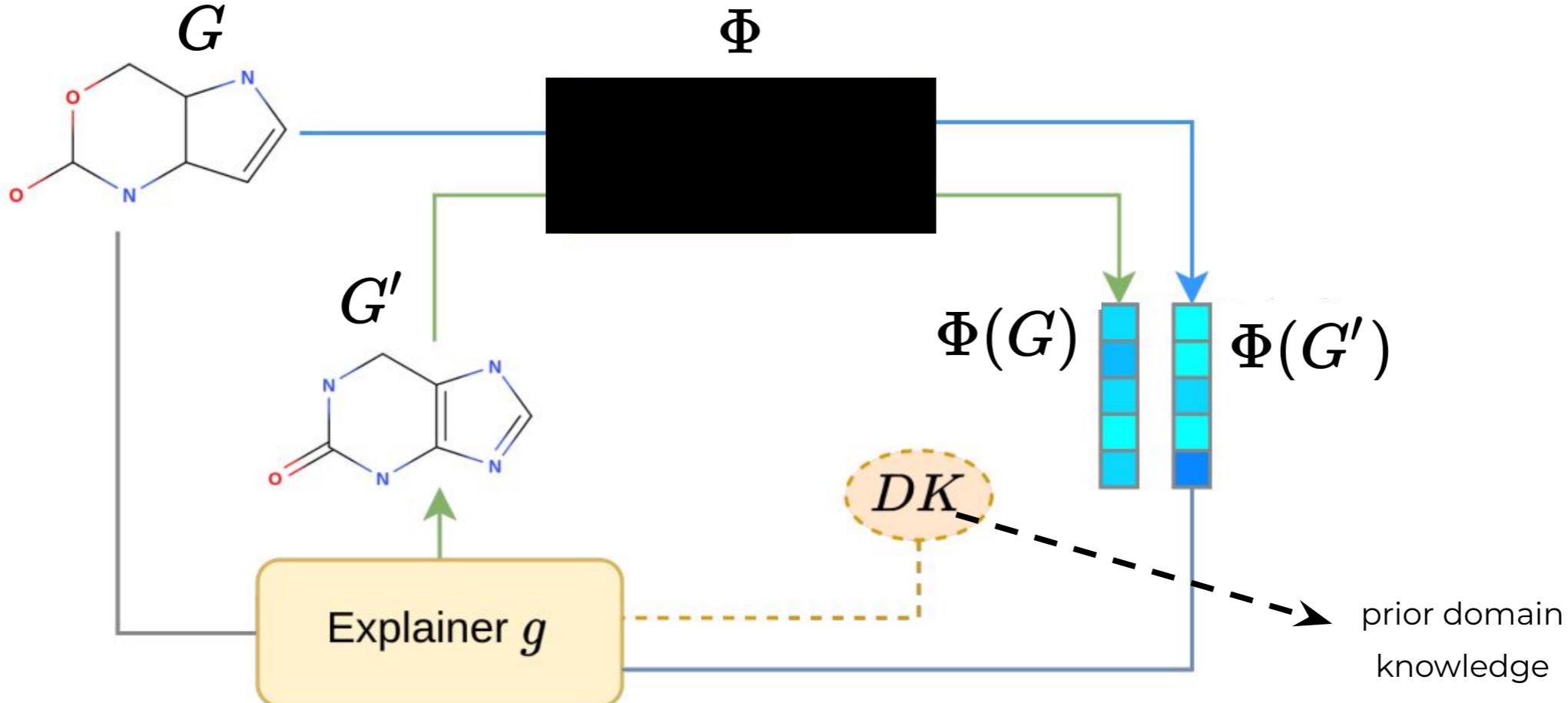
Wellawatte GP, Gandhi HA, Seshadri A, White AD. A Perspective on Explanations of Molecular Prediction Models. Journal of Chemical Theory and Computation. 2023 Mar 27;19(8):2149-60.

Tan, J., Geng, S., Fu, Z., Ge, Y., Xu, S., Li, Y. and Zhang, Y., 2022, April. Learning and evaluating graph neural network explanations based on counterfactual and factual reasoning. In Proceedings of the ACM Web Conference 2022 (pp. 1018-1027).

Ma, J., Guo, R., Mishra, S., Zhang, A. and Li, J., 2022. Clear: Generative counterfactual explanations on graphs. Advances in Neural Information Processing Systems, 35, pp.25895-25907.

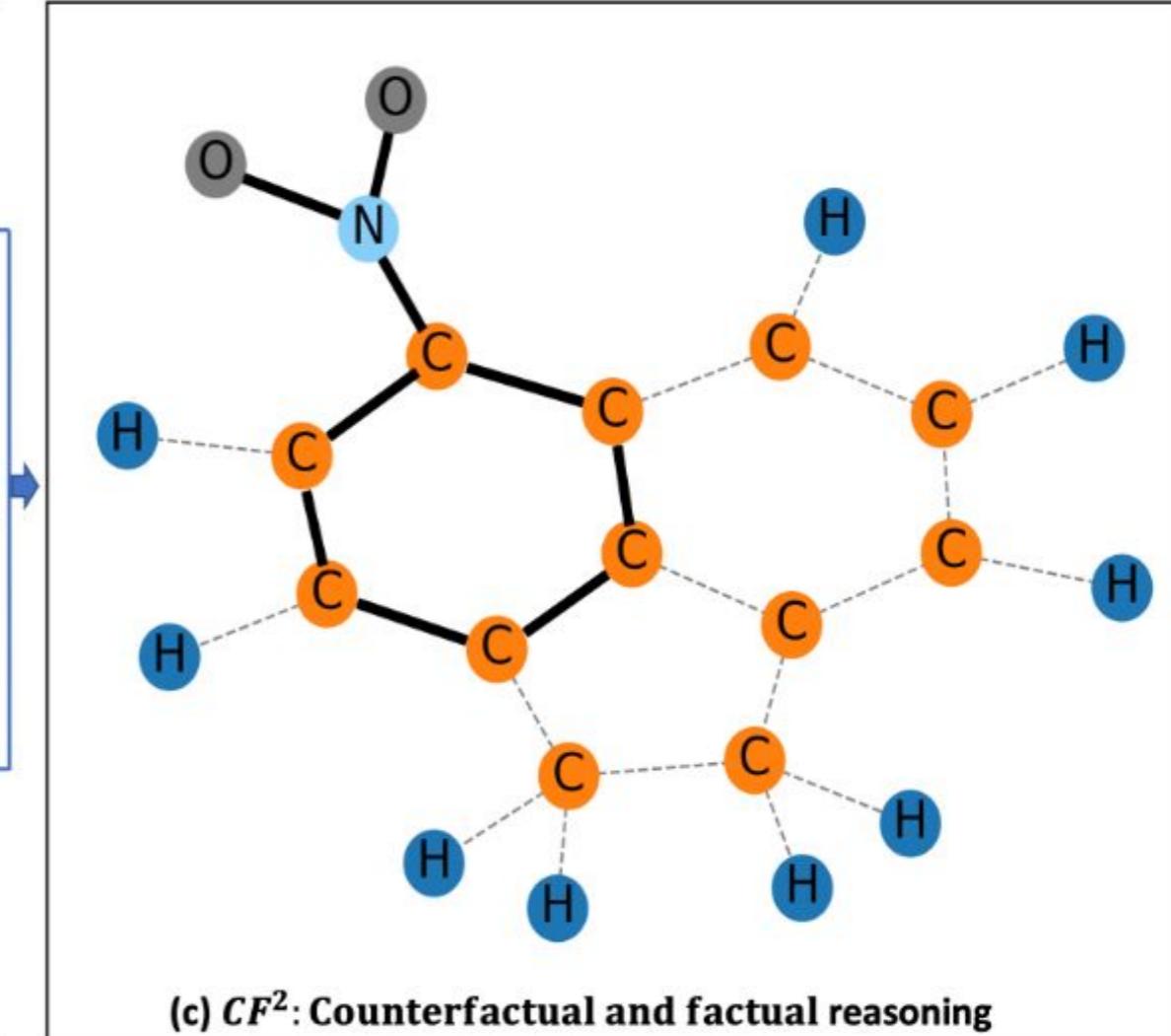
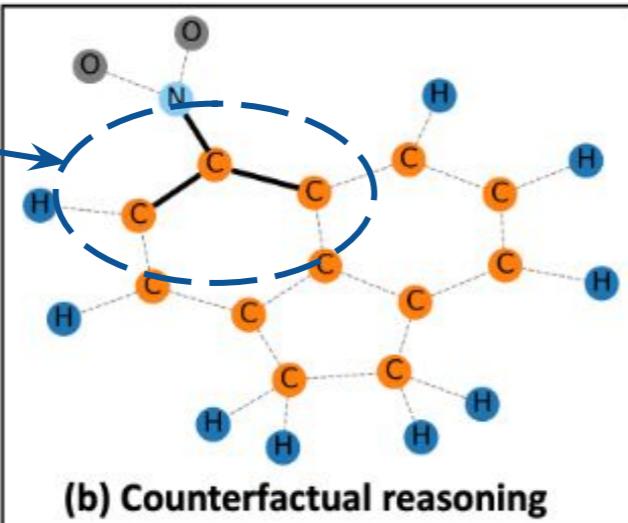
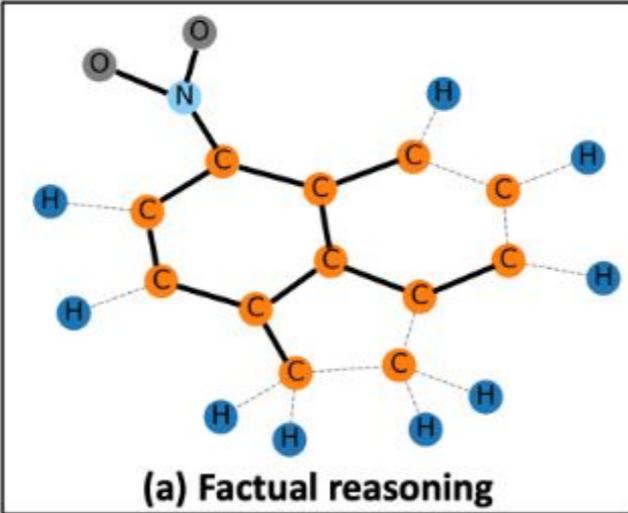
Prado-Romero MA, Prenkaj B, Stilo G. Robust Stochastic Graph Generator for Counterfactual Explanations. arXiv preprint arXiv:2312.11747. 2023 Dec 18.

# MEG (ARCHITECTURE)



# COUNTERFACTUAL AND FACTUAL ( $CF^2$ )

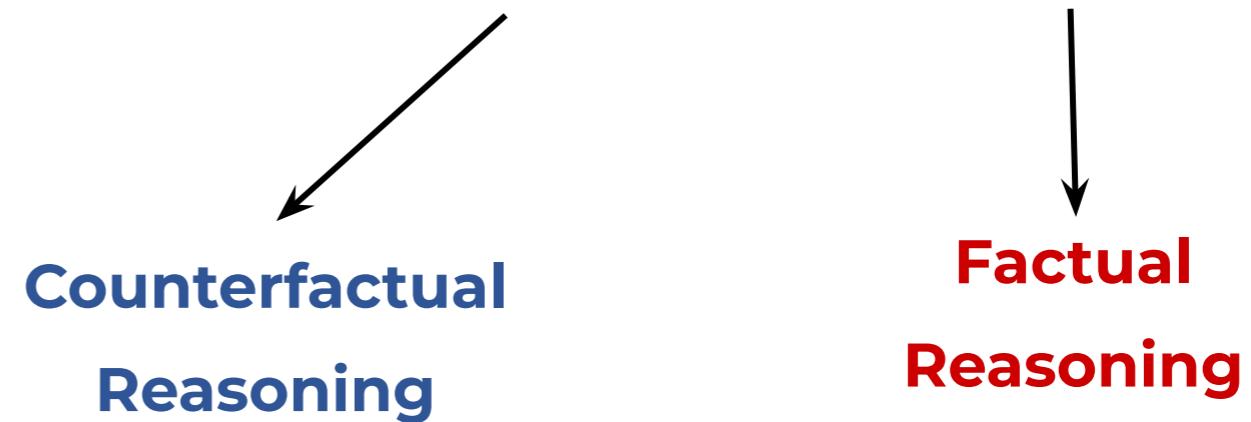
these edges, if  
**removed**,  
create the  
**counterfactual**



# CF<sup>2</sup> (NECESSITY & SUFFICIENCY)

---

An explanation needs to be **necessary** and **sufficient**



# CF<sup>2</sup> (SUFFICIENCY)

# Factual Reasoning

$$\Phi(A \ast M, X \ast F) = y$$

# Masked adjacency matrix

# Masked node feature vectors

Both masks are learned by the explainer



## Counterfactual Reasoning

$$\Phi(A - A*M, X - X*F) \neq y$$

Why do we subtract the masked adjacency matrix?

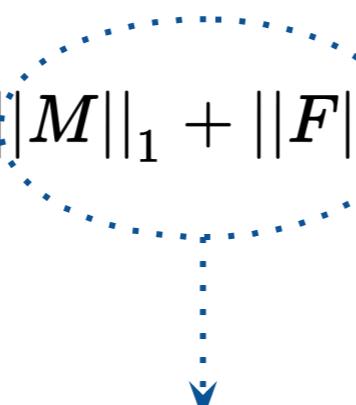


# CF<sup>2</sup> (RELAXED OPTIMIZATION)

---

$$L_f = \text{ReLU}(\gamma + P(\Phi(A * M, X * F) = \neg y) - S_f(M, F))$$

$$L_c = \text{ReLU}(\gamma - P(\Phi(A - A * M, X - X * F) = \neg y) - S_c(M, F))$$

$$\text{minimize } ||M||_1 + ||F||_1 + \lambda(\alpha L_f + (1 - \alpha)L_c)$$


**Complexity**

First work to treat **causality** when producing counterfactuals

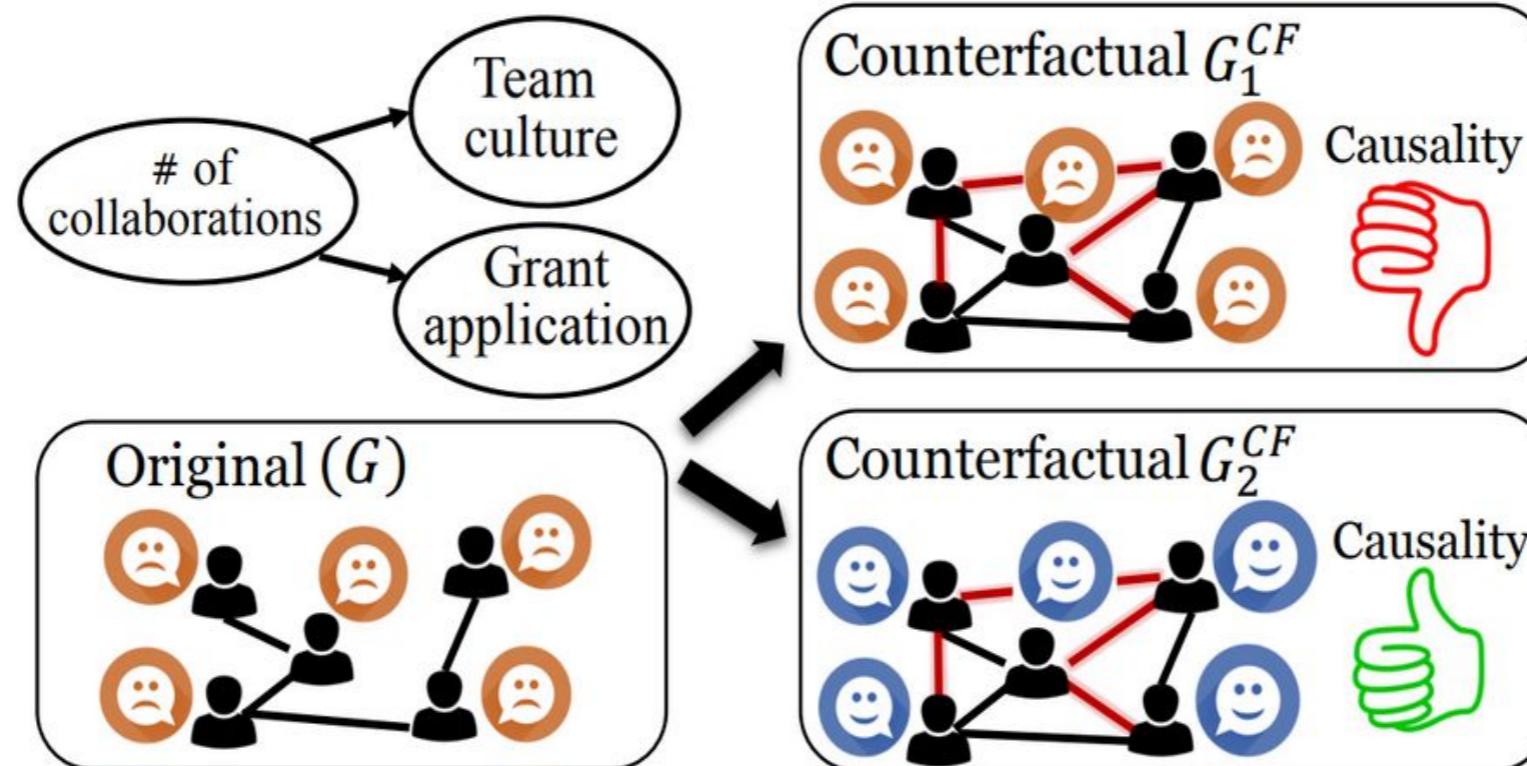
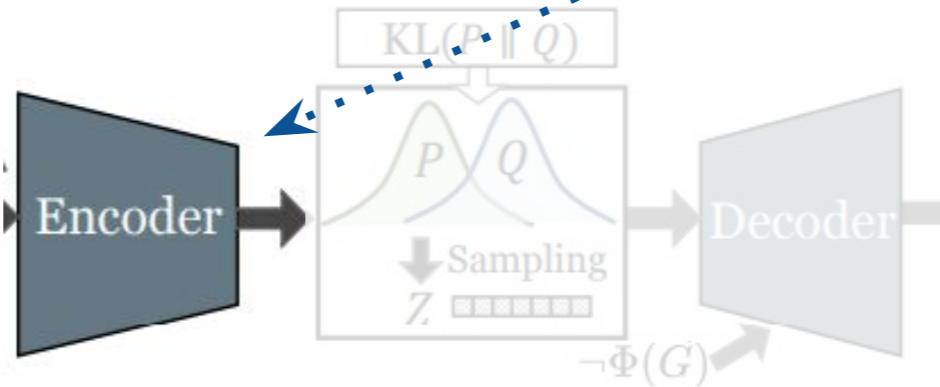


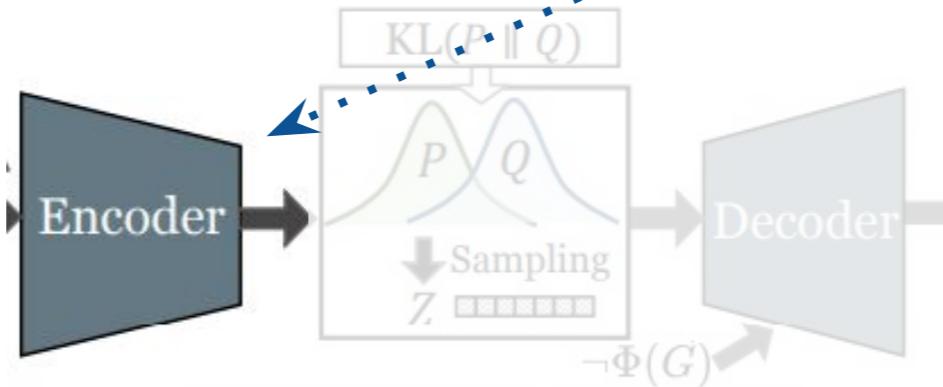
Image taken from: Ma J, Guo R, Mishra S, Zhang A, Li J. Clear: Generative counterfactual explanations on graphs. Advances in Neural Information Processing Systems. 2022 Dec 6;35:25895-907.

$$\mathcal{L} = \mathbb{E}_Q \left[ d(G, G') + \alpha \ell(\Phi(G'), \neg\Phi(G)) \right] + \text{KL}\left(Q(Z | G, \neg\Phi(G)) || P(Z | G, \neg\Phi(G))\right)$$

**Q is learned by the encoder**



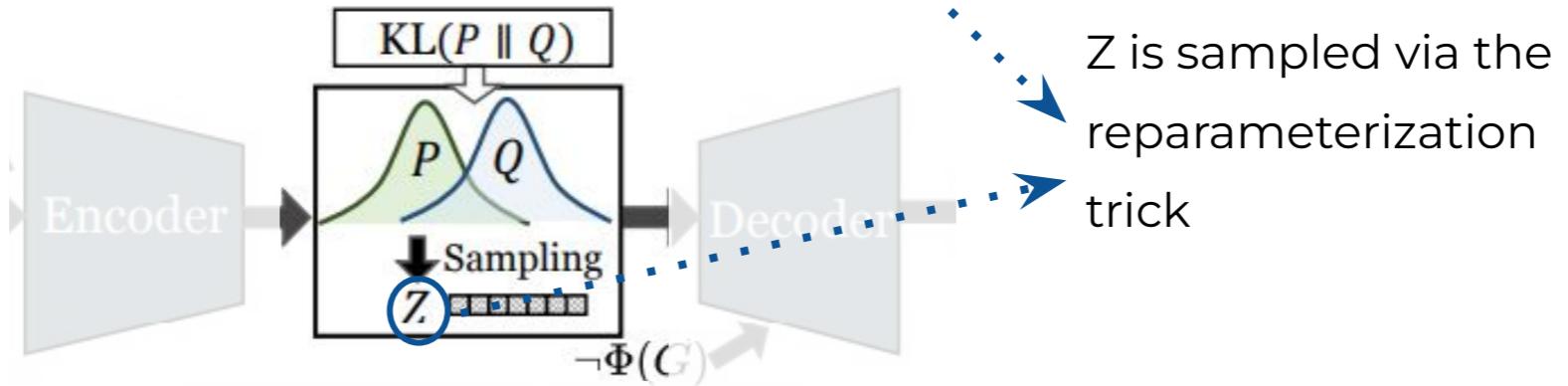
$$\mathcal{L} = \mathbb{E}_Q \left[ d(G, G') + \alpha \ell(\Phi(G'), \neg\Phi(G)) \right] + \text{KL}\left(Q(Z | G, \neg\Phi(G)) || P(Z | G, \neg\Phi(G))\right)$$



$$P(Z | G, \neg\Phi(G)) = \mathcal{N}\left(\mu_z(\neg\Phi(G)), \text{diag}(\sigma_z^2(\neg\Phi(G)))\right)$$

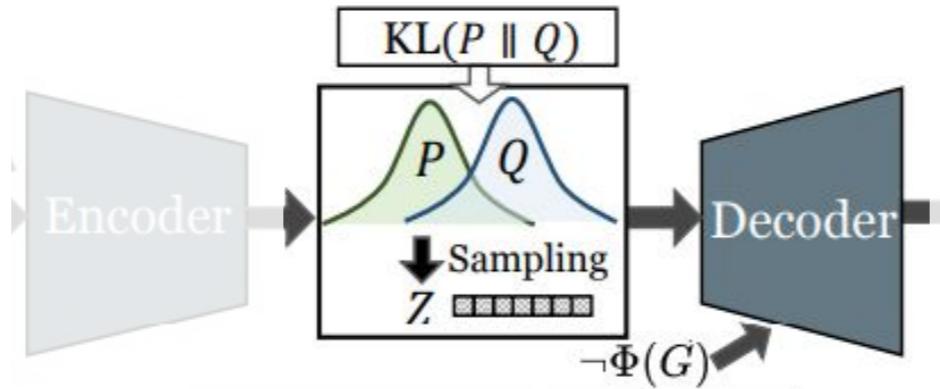
The Gaussian distribution as prior to enforce the learned distribution  $Q$  to be close to the prior by minimizing their KL divergence

$$\mathcal{L} = \mathbb{E}_Q \left[ d(G, G') + \alpha \ell(\Phi(G'), \neg\Phi(G)) \right] + \text{KL}\left(Q(\mathbf{Z}|G, \neg\Phi(G)) \parallel P(Z | G, \neg\Phi(G))\right)$$



Kingma DP, Welling M. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114. 2013 Dec 20.

- One can generate **multiple counterfactuals** from sampling multiple  $Z$
- The decoder produces a probabilistic graph where edges have weights  $\mathbb{R}_0^1$
- Binarize the graph according to the Bernoulli distribution



$$d(G, G') = d_A(A, \text{Bernoulli}(A')) + d_X(X, X')$$

distance between the original adjacency  
matrix and the generated one

distance between the original node  
features and the generated ones

- Cornerstone paper in the debate “Are generative counterfactual explanation approaches worth it?”
- **Besides CLEAR, all other explainers are discriminative**
- Uses Residual GANs to learn how to generate counterfactuals



# RSGG-CE (INTUITION)

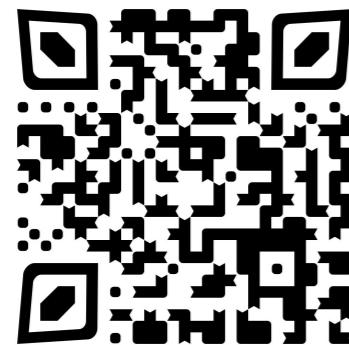
---

- Train the generator on graphs of class c (**the class to explain**)
- Train the discriminator on graphs different from class c and the synthetic data generated
- Because the generator needs to fool the discriminator, it'll learn to produce graphs of class **not c**
- Use the oracle to guide the generator in crossing the decision boundary



**Saturday 24th as Poster**  
on Hall A on the Exhibition Level

**Sunday 25th as Spotlight  
Paper** in Room 217 at 9:30



**CODE**

## **Robust Stochastic Graph Generator for Counterfactual Explanations**

**Mario Alfonso Prado-Romero<sup>\*1</sup>, Bardh Prenkaj<sup>\*2</sup>, Giovanni Stilo<sup>3</sup>**

<sup>1</sup> Gran Sasso Science Institute <sup>2</sup> Sapienza University of Rome <sup>3</sup> University of L'Aquila  
marioalfonso.prado@gssi.it, prenkaj@di.uniroma1.it, giovanni.stilo@univaq.it

<https://github.com/MarioTheOne/GRETEL>

### **Abstract**

Counterfactual Explanation (CE) techniques have garnered attention as a means to provide insights to the users engaging with AI systems. While extensively researched in domains such as medical imaging and autonomous vehicles, Graph Counterfactual Explanation (GCE) methods have been comparatively under-explored. GCEs generate a new graph similar to the original one, with a different outcome grounded on the underlying predictive model. Among these GCE techniques, those rooted in generative mechanisms have received relatively limited investigation despite demonstrating impressive accomplishments in other domains, such as artistic styles and natural language modelling. The preference for generative explainers stems from their capacity to generate counterfactual instances during inference, leveraging autonomously acquired perturbations of the input graph.

Prenkaj et al. 2021, 2020, 2023a; Verma, Mandal, and Gupta 2022; Wang, Yu, and Miao 2017).

Recently, deep learning (relying on GNNs (Scarselli et al. 2008)) has been beneficial in solving graph-based prediction tasks, such as community detection (Wu et al. 2022), link prediction (Wei et al. 2022), and session-based recommendations (Wu et al. 2019; Xu, Xi, and Wang 2021). Despite their remarkable performance, GNNs are black boxes, making them unsuitable for high-impact and high-risk scenarios. The literature has proposed several post-hoc explainability methods to understand *what is happening under the hood* of the prediction models. Specifically, counterfactual explainability is useful to understand how modifications in the input lead to different outcomes. Similarly, a recent field in Graph Counterfactual Explainability (GCE) has emerged





# DATASETS AND EVALUATION METRICS

# GRAPH CLASSIFICATION DATASETS IN THE WILD

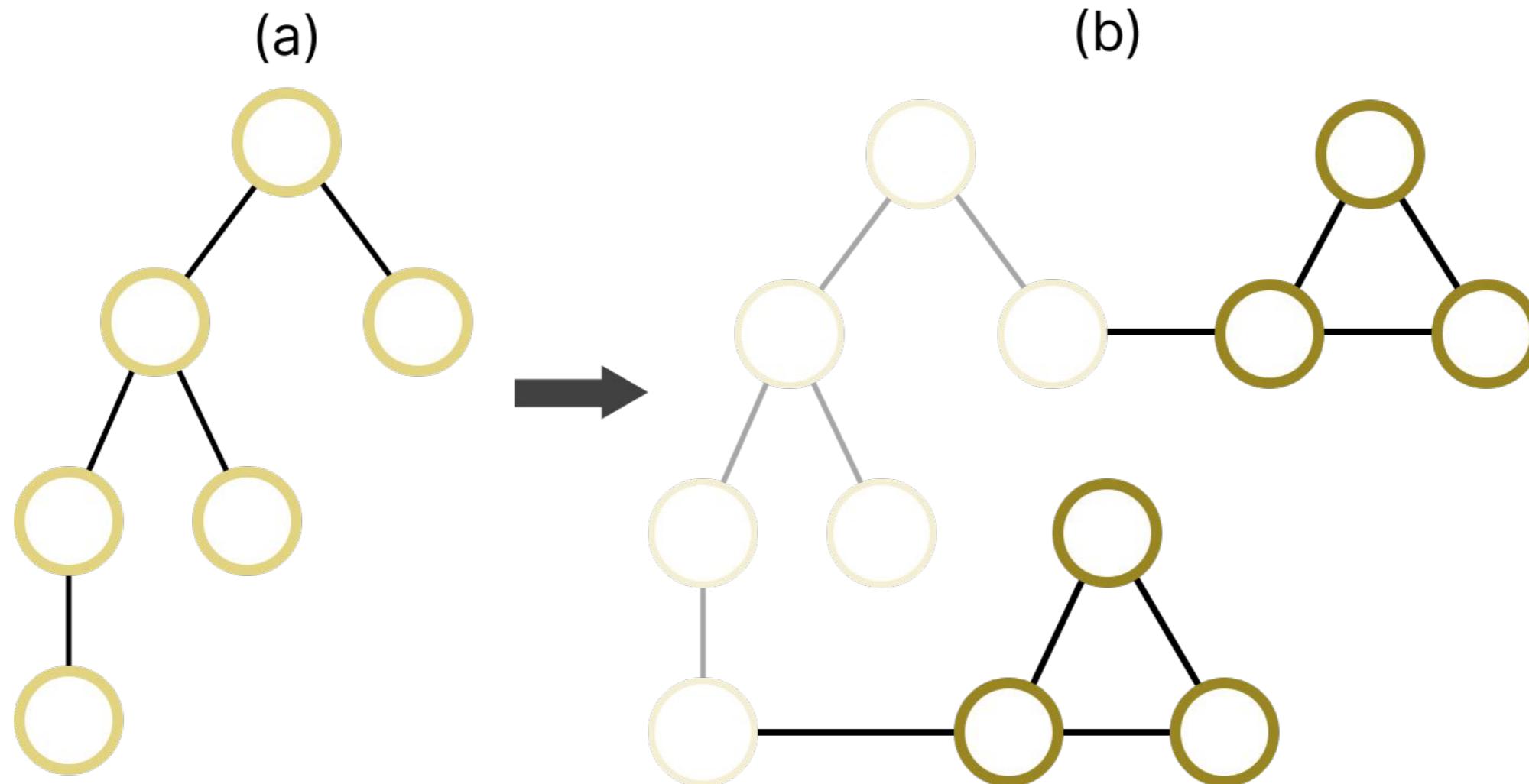
Dataset	Domain	Publicly Available Repository (Data or Code)	Used by
Tree-Cycles [98]	synthetic	<a href="https://github.com/RexYing/gnn-model-explainer">https://github.com/RexYing/gnn-model-explainer</a>	[6, 15, 43, 81]
Tree-Grid [98]	synthetic	<a href="https://github.com/RexYing/gnn-model-explainer">https://github.com/RexYing/gnn-model-explainer</a>	[6, 15, 43]
Tree-Infinity	synthetic	<a href="https://github.com/MarioTheOne/GRETEL">https://github.com/MarioTheOne/GRETEL</a>	[68]
BA-Shapes [98]	synthetic	<a href="https://github.com/RexYing/gnn-model-explainer">https://github.com/RexYing/gnn-model-explainer</a>	[6, 15, 43, 81]
BA-Community [98]	synthetic	<a href="https://github.com/RexYing/gnn-model-explainer">https://github.com/RexYing/gnn-model-explainer</a>	[6]
BA-2motifs [44]	synthetic	<a href="https://github.com/flyingdoog/PGEExplainer">https://github.com/flyingdoog/PGEExplainer</a>	[6, 81]
ADHD [13]	-omics	<a href="https://github.com/MarioTheOne/GRETEL/tree/main/data/datasets/adhd">https://github.com/MarioTheOne/GRETEL/tree/main/data/datasets/adhd</a>	[1]
ASD [19, 37]	-omics	<a href="https://github.com/MarioTheOne/GRETEL/tree/main/data/datasets/autism/asd">https://github.com/MarioTheOne/GRETEL/tree/main/data/datasets/autism/asd</a>	[1]
BBBP [50]	molecular	<a href="https://www.kaggle.com/datasets/mmmelahi/cheminformatics?select=bbbp.zip">https://www.kaggle.com/datasets/mmmelahi/cheminformatics?select=bbbp.zip</a>	[91]
HIV [20, 24, 70]	molecular	<a href="https://www.kaggle.com/datasets/mmmelahi/cheminformatics?select=hiv.zip">https://www.kaggle.com/datasets/mmmelahi/cheminformatics?select=hiv.zip</a>	[32, 91]
Ogbg-molhiv [31]	molecular	<a href="https://huggingface.co/datasets/OGB/ogbg-molhiv">https://huggingface.co/datasets/OGB/ogbg-molhiv</a>	[45]
Mutagenicity [34]	molecular	<a href="https://ls11-www.cs.tu-dortmund.de/people/morris/graphkerneldatasets/Mutagenicity.zip">https://ls11-www.cs.tu-dortmund.de/people/morris/graphkerneldatasets/Mutagenicity.zip</a>	[6, 32, 81]
NCI1 [87]	molecular	<a href="https://ls11-www.cs.tu-dortmund.de/people/morris/graphkerneldatasets/NCI1.zip">https://ls11-www.cs.tu-dortmund.de/people/morris/graphkerneldatasets/NCI1.zip</a>	[6, 32, 81]
TOX21 [35]	molecular	<a href="https://tripod.nih.gov/tox21/challenge/data.jsp">https://tripod.nih.gov/tox21/challenge/data.jsp</a>	[58]
ESOL [95]	molecular	<a href="https://github.com/deepchem/deepchem">https://github.com/deepchem/deepchem</a>	[58, 81]
Proteins [11]	molecular	<a href="https://chrsmrrs.github.io/datasets/docs/datasets/">https://chrsmrrs.github.io/datasets/docs/datasets/</a>	[32]
Davis [21]	molecular	<a href="http://staff.cs.utu.fi/~aatapa/data/DrugTarget/">http://staff.cs.utu.fi/~aatapa/data/DrugTarget/</a>	[55]
PDBBind [89]	molecular	<a href="http://www.pdbbind.org.cn/">http://www.pdbbind.org.cn/</a>	[55]
CiteSeer [26]	social	<a href="https://linqs.org/datasets/">https://linqs.org/datasets/</a>	[40, 81]
IMDB-M [96]	social	<a href="https://virginia.app.box.com/s/941v9pwh83lfw5vnwfbgcrtlsoivg5j">https://virginia.app.box.com/s/941v9pwh83lfw5vnwfbgcrtlsoivg5j</a>	[45]
CORA [51]	social	<a href="https://relational.fit.cvut.cz/dataset/CORA">https://relational.fit.cvut.cz/dataset/CORA</a>	[40]
Musae-Facebook [71]	social	<a href="https://www.kaggle.com/datasets/rozemberczki/musae-facebook-pagepage-network">https://www.kaggle.com/datasets/rozemberczki/musae-facebook-pagepage-network</a>	[40]
LastFM [72]	social	<a href="https://github.com/gusye1234/LightGCN-PyTorch/tree/master/data/lastfm">https://github.com/gusye1234/LightGCN-PyTorch/tree/master/data/lastfm</a>	[16]
Yelp [90]	social	<a href="https://github.com/gusye1234/LightGCN-PyTorch/tree/master/data/yelp2018/">https://github.com/gusye1234/LightGCN-PyTorch/tree/master/data/yelp2018/</a>	[16]

Prado-Romero MA, Prenkaj B, Stilo G, Giannotti F. A survey on graph counterfactual explanations: definitions, methods, evaluation. arXiv preprint arXiv:2210.12089. 2022 Oct 21.



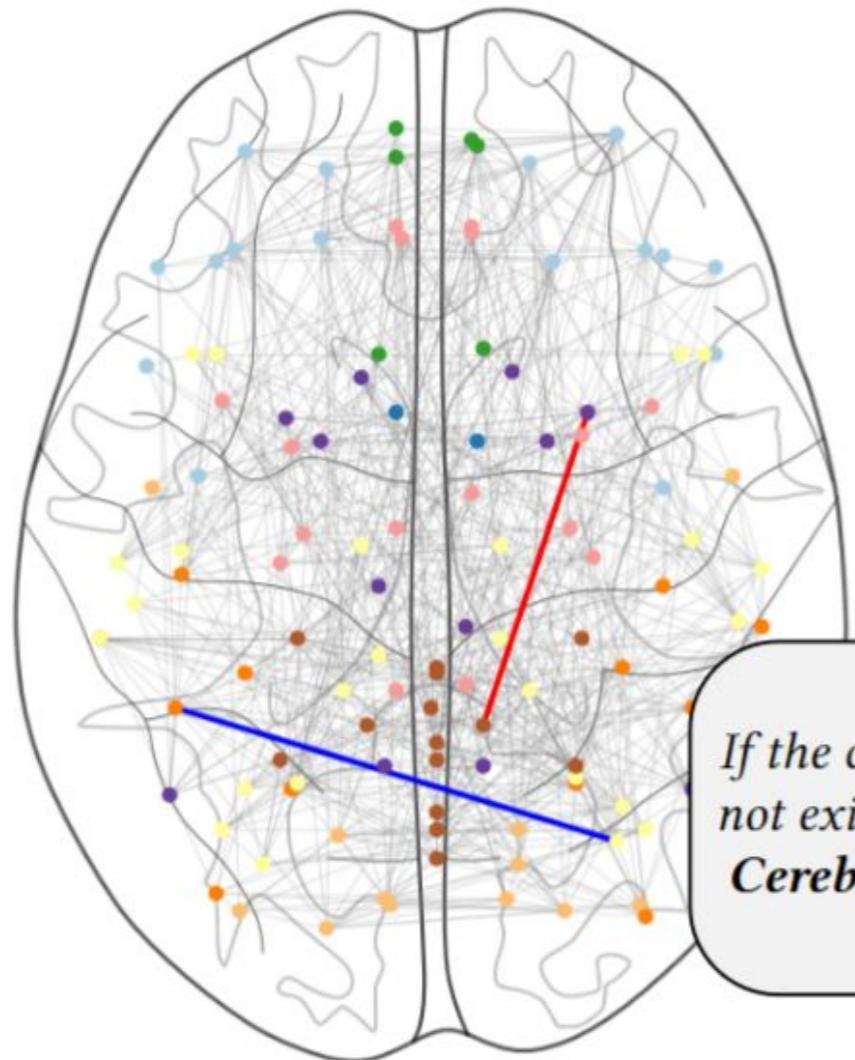
# TREECYCLES

---



Ying Z, Bourgeois D, You J, Zitnik M, Leskovec J. Gnnexplainer: Generating explanations for graph neural networks. Advances in neural information processing systems. 2019;32.

# AUTISM SPECTRUM DISORDER (ASD)

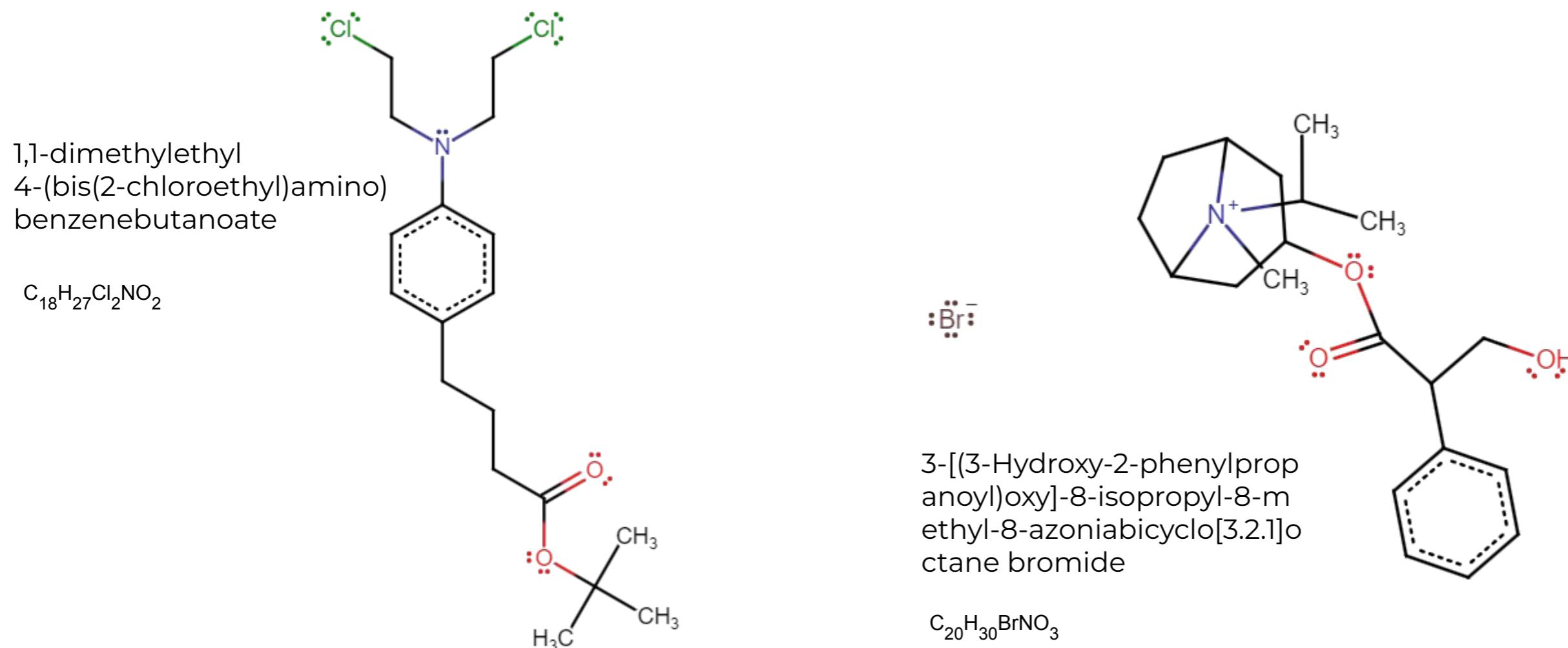


Frontal Lobe	Deep Grey Matter
Insula	Cerebellum
Cingulate	Others
Temporal Lobe	C-Edges (+)
Occipital Lobe	C-Edges (-)
Parietal Lobe	O-Edges

*Patient USM\_0050453 is classified as ASD (class 1). If the connection between **Putamen\_R** and **Cerebellum\_9\_R** did not exist and instead the connection between **Parietal\_Inf\_L** and **Cerebellum\_Crus2\_R** existed, then the patient would have been classified as TD (class 0).*

Abbate C, Bonchi F. Counterfactual graphs for explainable classification of brain networks. In Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining 2021 Aug 14 (pp. 2495-2504).

# BLOOD BARRIER PENETRATION PREDICTION (BBBP)



Martins IF, Teixeira AL, Pinheiro L, Falcao AO. A Bayesian approach to in silico blood-brain barrier penetration modeling. Journal of chemical information and modeling. 2012 Jun 25;52(6):1686-97.

# EVALUATION METRICS

---

**Correctness** (Validity) 

$$\Omega(G, G') = \mathbf{1}[\Phi(G) \neq \Phi(G')]$$

We're answering the question: “*Can the explainer produce an instance with a **different classification** from the original instance?*”



# EVALUATION METRICS

---

## Graph Edit Distance (GED)

$$\text{GED}(G, G') = \min_{\{p_1, \dots, p_n\} \in \mathcal{P}(G, G')} \sum_{i=1}^n \omega(p_i)$$

We're answering the question: “How **similar** is the produced counterfactual candidate w.r.t. original input?”



# GED TRICK

---

$$\text{GED}(G, G') = \min_{\{p_1, \dots, p_n\} \in \mathcal{P}(G, G')} \sum_{i=1}^n \omega(p_i)$$


$$\sum_i \sum_j |A_{i,j} - A'_{i,j}|$$

Just sum the absolute element-wise difference to get the edges that were changed in the counterfactual w.r.t. the original input

*(Works only if we have binary adjacency matrices, and their dimensionalities remain the same)*

# EVALUATION METRICS

---

## Fidelity

$$\Psi(G, G') = \mathbf{1}[\Phi(G) = y_G] - \mathbf{1}[\Phi(G') = y_G]$$

We're answering the question: “*How **faithful** the explanations are to the oracle considering their correctness?*”



# LET'S DE-RECONSTRUCT FIDELITY

---

Fidelity 

$$\Psi(G, G') = \mathbf{1}[\Phi(G) = y_G] - \mathbf{1}[\Phi(G') = y_G]$$

We want the oracle to be correct in predicting the ground truth of the original instance (**aka accuracy**)



# LET'S DE-RECONSTRUCT FIDELITY

---

Fidelity 

$$\Psi(G, G') = \mathbf{1}[\Phi(G) = y_G] - \mathbf{1}[\Phi(G') = y_G]$$

We don't want the oracle to predict the same class as the ground truth



# LET'S DE-RECONSTRUCT FIDELITY

---

## Fidelity can have 3 values

- **+1** → both the explainer and oracle are working correctly
- **0 & -1** → something is wrong with the explainer or the oracle

*When the oracle has perfect accuracy, fidelity is equal to correctness*





# TAKEAWAY LESSONS

# SOA PERFORMANCE

**Accuracy = 100%**

Dataset	Method	GED ↓	Correctness ↑	Fidelity ↑
Tree-Cycles	RAND@5	$92.18 \pm 5.44$	$0.55 \pm 0.50$	$0.55 \pm 0.50$
	RAND@10	$123.74 \pm 7.43$	$0.51 \pm 0.50$	$0.51 \pm 0.50$
	RAND@15	$147.93 \pm 8.26$	$0.58 \pm 0.50$	$0.58 \pm 0.50$
	DCE	$50.36 \pm 0.00$	$1.00 \pm 0.00$	$1.00 \pm 0.00$
	OBS	$57.31 \pm 0.03$	$0.96 \pm 0.01$	$0.96 \pm 0.01$
	DDBS	$71.79 \pm 0.24$	$0.59 \pm 0.01$	$0.59 \pm 0.01$
	MACCS	—	—	—
	CLEAR	$79.76 \pm 3.60$	$0.53 \pm 0.10$	$0.53 \pm 0.10$
	CF <sup>2</sup>	$31.54 \pm 0.12$	$0.47 \pm 0.10$	$0.47 \pm 0.10$
	MEG	$159.70 \pm 1.34$	$0.53 \pm 0.00$	$0.53 \pm 0.00$

Naïve baseline is better than intensive and complicated generative approach

Chance level correctness!



# SOA PERFORMANCE

**Accuracy = 79%**

Dataset	Method	GED ↓	Correctness ↑	Fidelity ↑
ASD	RAND@5	$618.06 \pm 8.27$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	RAND@10	$1152.93 \pm 20.19$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	RAND@15	$1600.78 \pm 18.22$	$0.00 \pm 0.00$	$0.00 \pm 0.00$
	DCE	$1011.69 \pm 0.00$	$1.00 \pm 0.00$	$0.54 \pm 0.00$
	OBS	$9.89 \pm 0.11$	$1.00 \pm 0.00$	$0.54 \pm 0.00$
	DDBS	$11.79 \pm 0.29$	$1.00 \pm 0.00$	$0.54 \pm 0.00$
	MACCS	—	—	—
	CLEAR	$1739.60 \pm 131.16$	$0.47 \pm 0.13$	$0.25 \pm 0.18$
	CF <sup>2</sup>	$655.49 \pm 2.87$	$0.46 \pm 0.09$	$0.37 \pm 0.15$
	MEG	×	×	×

Generative approaches don't even compare to the baselines



# SOA PERFORMANCE

**Accuracy = 86%**

Dataset	Method	GED ↓	Correctness ↑	Fidelity ↑
	RAND@5	$30.98 \pm 33.27$	$0.85 \pm 0.35$	$0.62 \pm 0.69$
	RAND@10	$52.98 \pm 58.96$	$0.86 \pm 0.35$	$0.65 \pm 0.66$
	RAND@15	$82.97 \pm 137.37$	$0.85 \pm 0.36$	$0.61 \pm 0.69$
	DCE	$27.92 \pm 0.12$	$1.00 \pm 0.00$	$0.72 \pm 0.00$
	OBS	$0.00 \pm 0.00$	$0.00 \pm 0.00$	$0.61 \pm 0.00$
BBBP	DDBS	×	×	×
	MACCS	$11.23 \pm 0.08$	$0.40 \pm 0.00$	$0.23 \pm 0.00$
	CLEAR@1	$27056.29 \pm 9.69$	$0.87 \pm 0.02$	$0.64 \pm 0.03$
	CLEAR@5	$26711.57 \pm 112.67$	$0.85 \pm 0.02$	$0.62 \pm 0.03$
	CLEAR@15	$25986.66 \pm 170.41$	$0.85 \pm 0.01$	$0.62 \pm 0.06$
	CF <sup>2</sup>	$25.72 \pm 0.63$	$0.85 \pm 0.02$	$0.63 \pm 0.03$
	MEG	$269.35 \pm 0.39$	$0.51 \pm 0.04$	$0.32 \pm 0.04$

Removing the factual graph  
can get you good results...

Good correctness but too far  
away from the input



# OUR CRITIQUES (SOME)

---

- Most datasets are toy-like, which don't allow for GNNs (oracles) to be correctly trained
  - See Abrate & Bonchi's paper: *customly written oracle which looks like a linear separator with fixed slope and y-intercept*
- To date, only RSGG-CE can compare against search-based explainers.  
*Why aren't the other works comparing to them as baselines!?*
- Let's define a standard evaluation benchmark: correctness, *GED*, and *fidelity must be included in all future proposals*





# WHAT'S NEXT?

# OPEN RESEARCH QUESTIONS

---

- Are generative counterfactual explainers worth it?
- What's the difference between **counterfactual explanations** and **adversarial attacks**?
- How sure are we about the generated counterfactuals? Can we incorporate **uncertainty** in them?

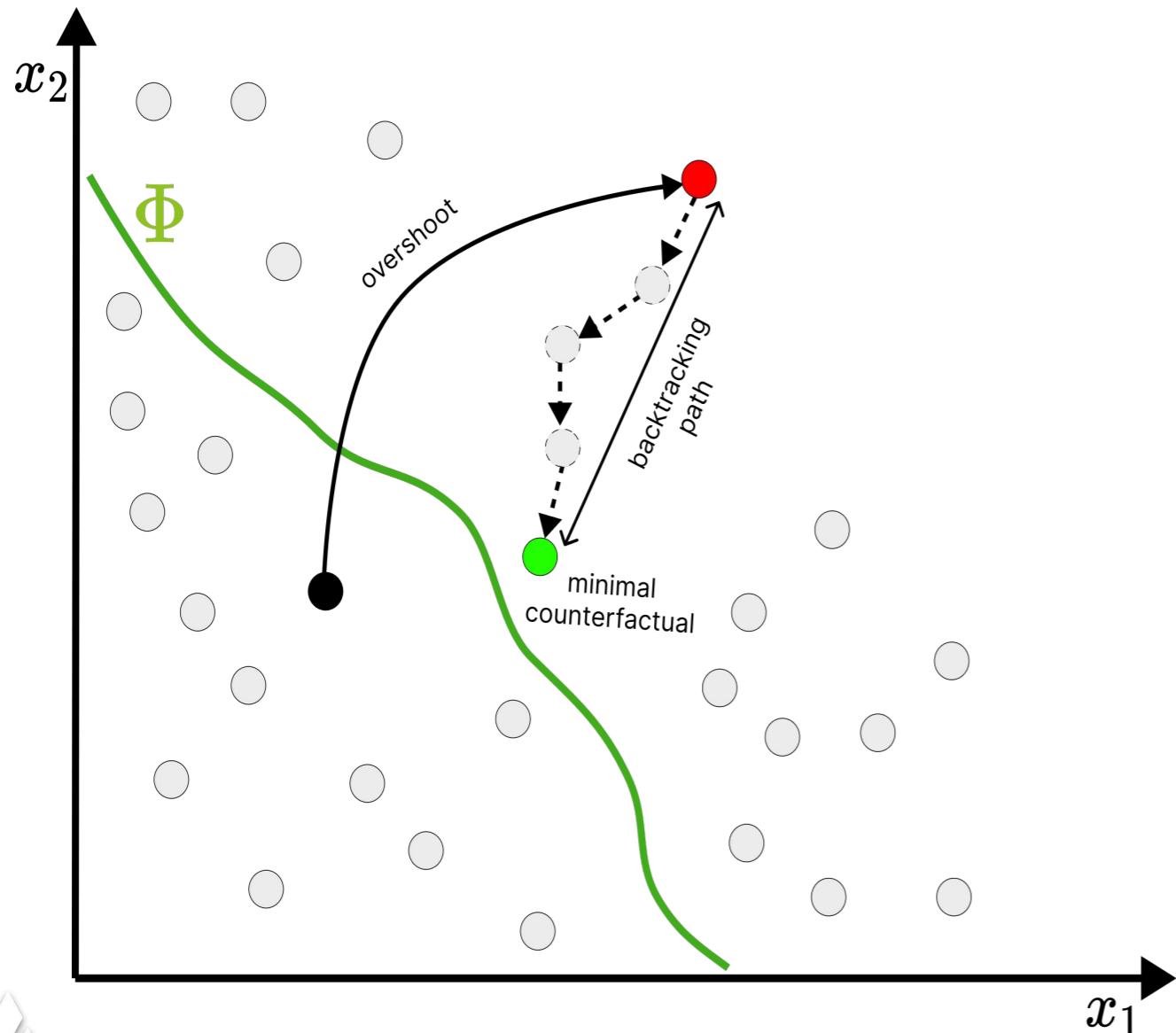


# OPEN RESEARCH QUESTIONS

---

- Are the produced **counterfactuals actionable?** How to incorporate **domain knowledge** into the explanation methods?
- How to ensure our explanations methods are **stable** and **robust**, producing **similar explanations for similar instances?**
- How can we **backtrack near the decision boundary** once we overshoot on the other side **to produce minimal counterfactuals?**  
(see next slide)

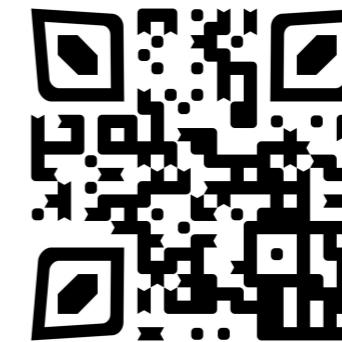




Can we learn the backtracking  
steps one-by-one?

Looks like an RL problem now...

**Care to help us? Work with us  
then**



# Thanks for your attention!



## QUESTIONS?



Mario A. Prado-Romero  
[marioalfonso.prado@gssi.it](mailto:marioalfonso.prado@gssi.it)



Bardh Prenkaj  
[prenkaj@di.uniroma1.it](mailto:prenkaj@di.uniroma1.it)

based on ACM survey



<https://dl.acm.org/doi/abs/10.1145/3618105>



Giovanni Stilo  
[giovanni.stilo@univaq.it](mailto:giovanni.stilo@univaq.it)