# Opportunities and Challenges in GraphML for Biomedicine
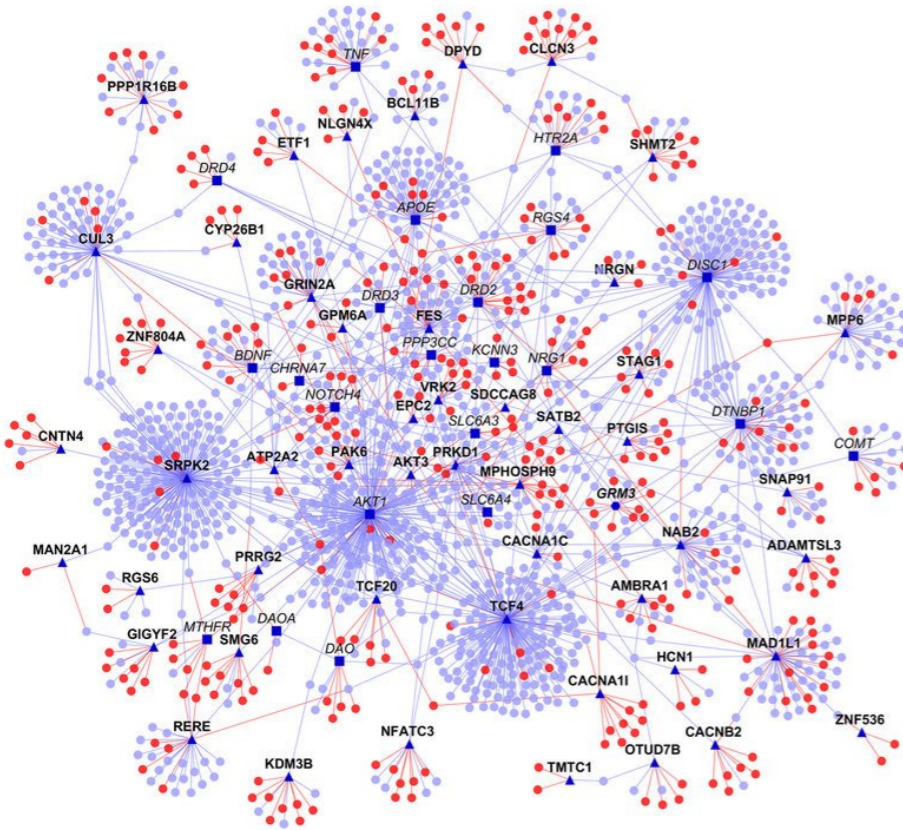
Megha Khosla
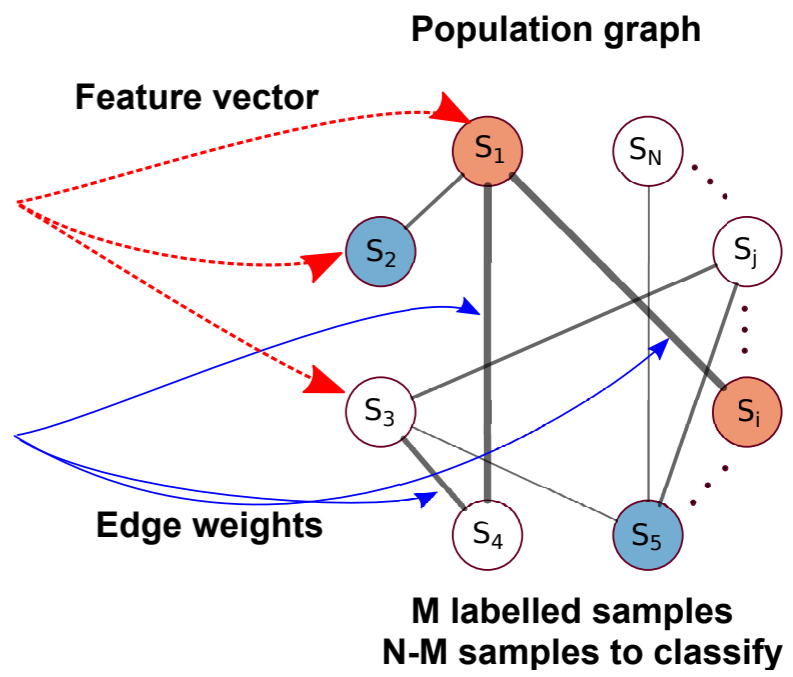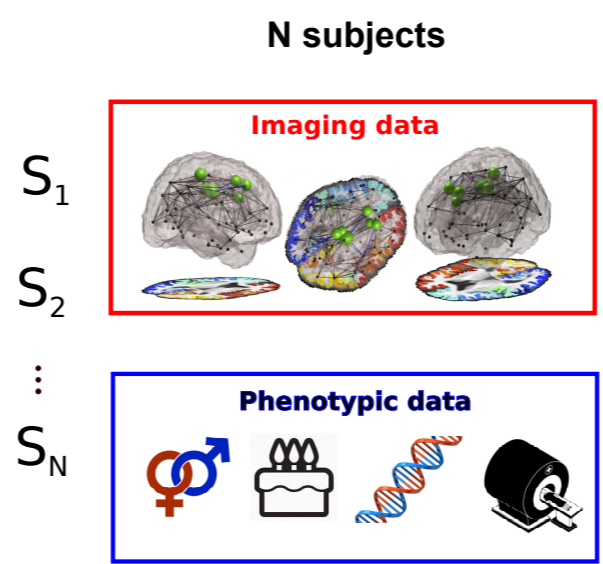
https://khosla.github.io
m.khosla@tudelft.nl

**TU**Delft

# Graphs in biomedicine



**Protein interaction network**

Image Source : wikipedia



N subjects

Imaging data

$S_1$

$S_2$

$S_N$

Phenotypic data

Population graph

Feature vector

Edge weights

$S_1$ $S_N$ $S_2$ $S_j$ $S_3$ $S_i$ $S_4$ $S_5$

M labelled samples
N-M samples to classify

**Patient Network**

Image Source : Parisot et al.

# Graph Machine Learning (GraphML)



Look-up table

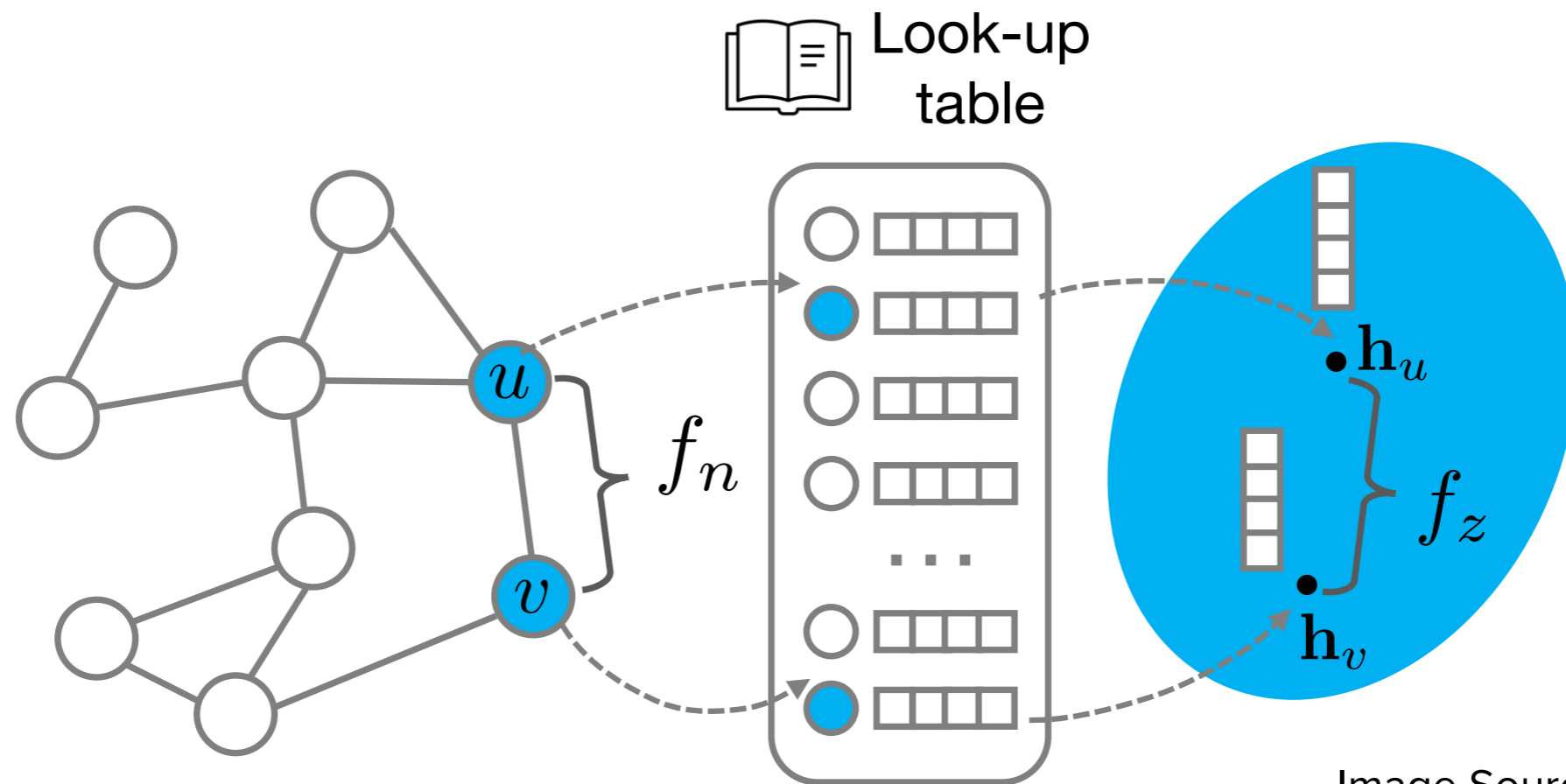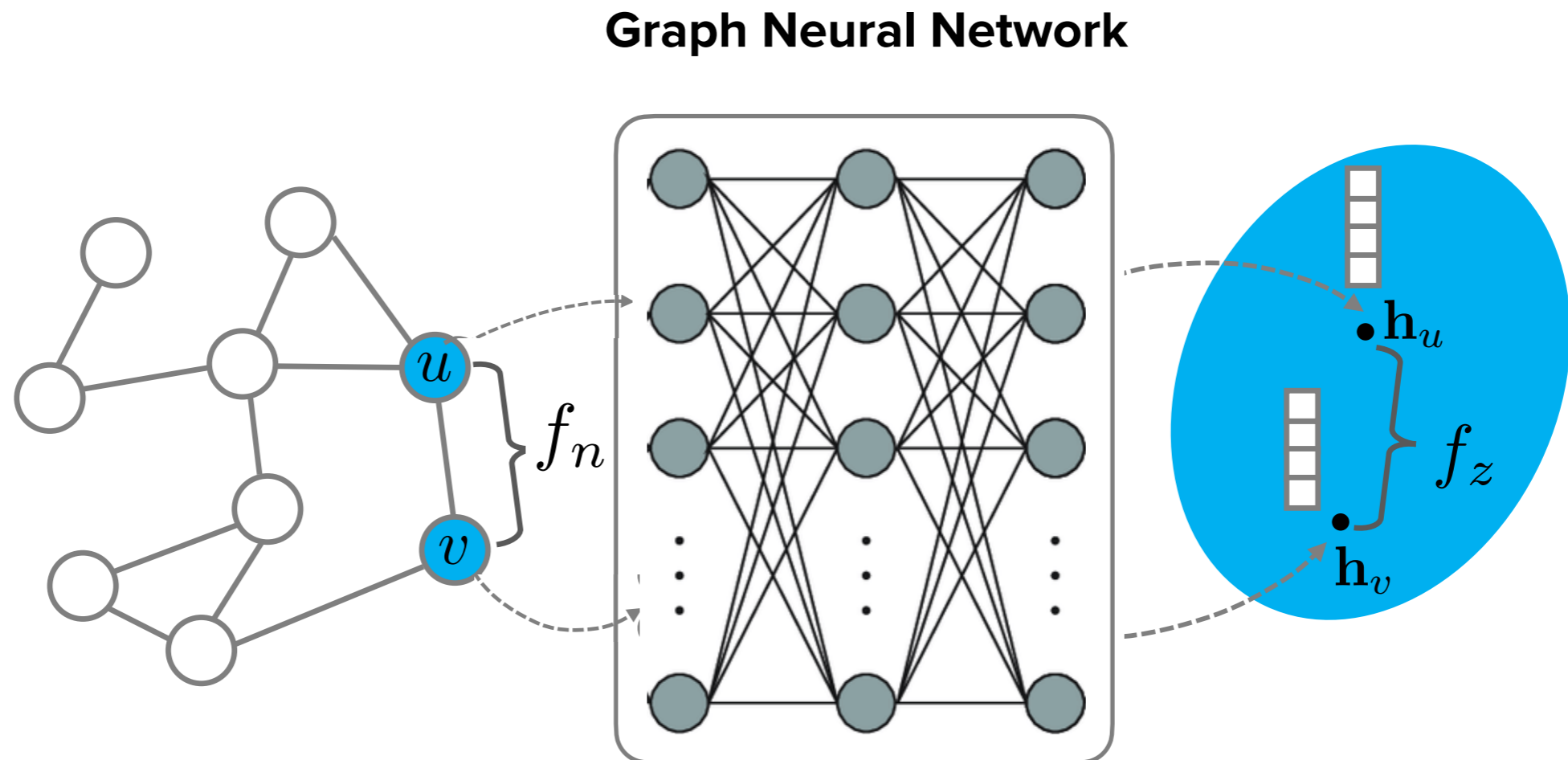$f_n$

$f_z$

$\mathbf{h}_u$

$\mathbf{h}_v$

$L_2$

$L_1$

Image Source: [Li et al., 2022]

**Shallow Network Embedding Methods**

Examples : **DeepWalk, Node2Vec, NERD, HOPE**

3

# Graph Machine Learning

**Graph Neural Network**



Examples :

**GCN, GAT, GIN**

# Applications of GraphML in Biomedicine

**Biological problems**

- Predict new human-pathogen protein interactions

- Predict new miRNA-disease associations

**Main Challenges**
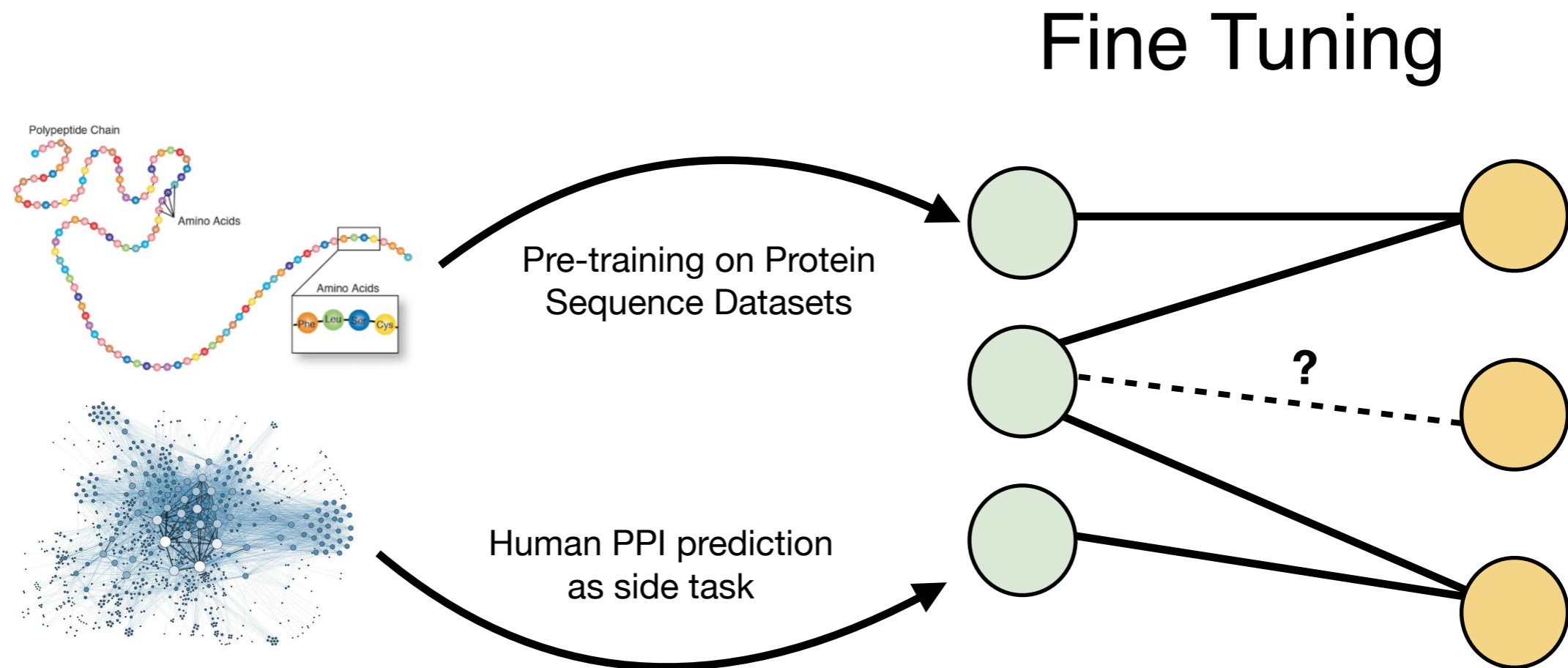
- Data scarcity

- Data Bias

**Other common issues**

- Wrong evaluation setups leading to data leakage

- Limited and biased train-test data

- N. Dong, J. Schrader, S. Mücke, M. Khosla, "*A Message Passing framework with Multiple data integration for miRNA-Disease association prediction*", In **Scientific Reports**, 2022.
- N. Dong, S. Mücke, M. Khosla, "*MuCoMiD: A Multitask graph Convolutional Learning Framework for miRNA-Disease Association Prediction*", in IEEE/ACM **Transactions on Computational Biology and Bioinformatics** 2022
- N. Dong, G. Brogden, G. Gerold, M. Khosla, "A multi-task transfer learning framework for the prediction of virus-human protein-protein interactions", **BMC Bioinformatics**, 2021.
- N.Dong, M.Khosla, Towards a consistent evaluation of miRNA-disease association prediction models. In IEEE International Conference on Bioinformatics and Biomedicine (**BIBM**), 2020
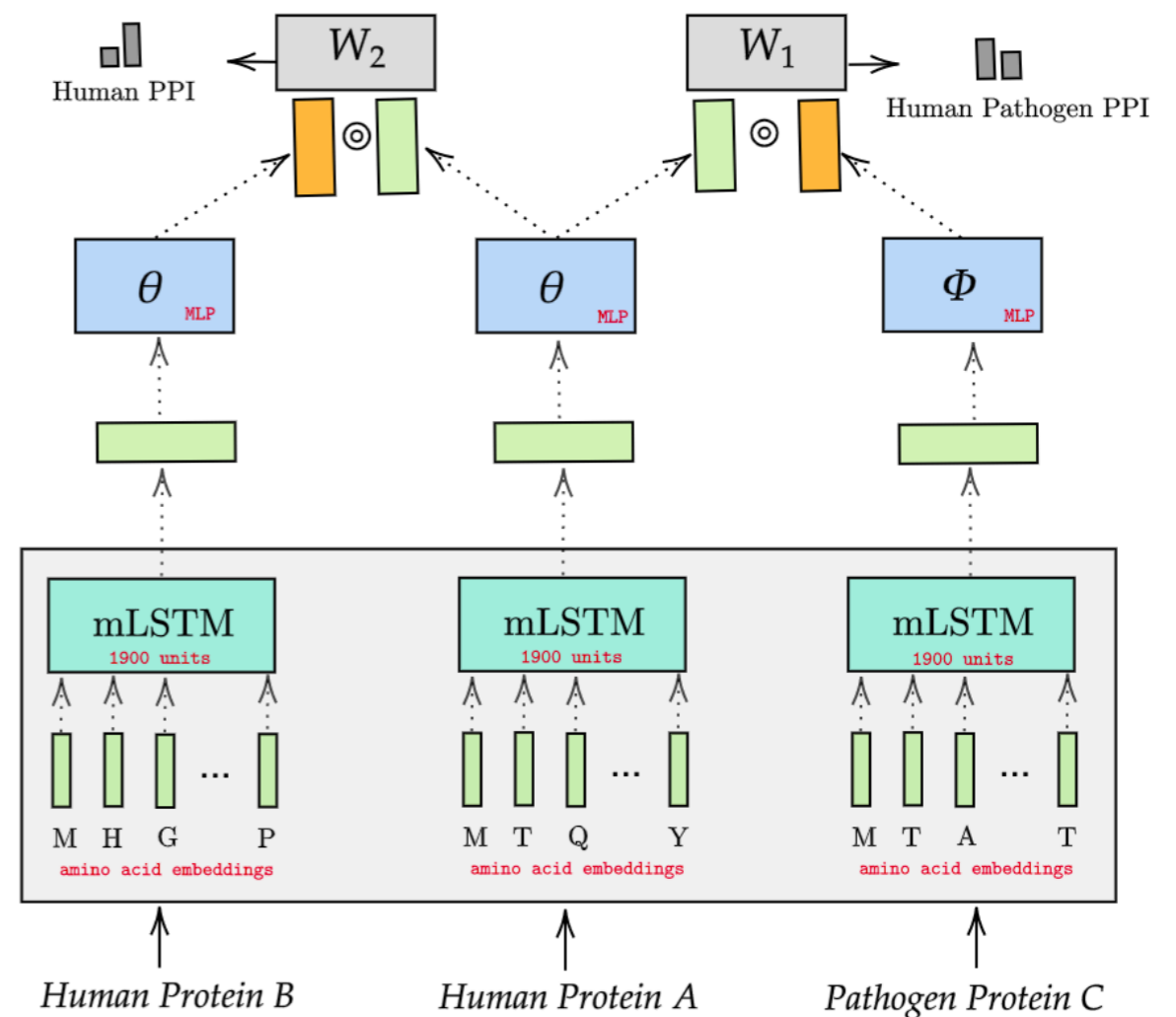
# Predicting protein-pathogen protein interactions

How to use inductive biases from multiple sources of information to overcome challenges of learning under low data regimes?



Fine Tuning

Pre-training on Protein Sequence Datasets

Human PPI prediction as side task

# Join learning framework

- Powerful input protein representations learnt over 24 million protein sequences

- Multitask learning framework using graph reconstruction losses

- Besides strong results on public datasets we could identify COVID 19 top receptor



https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-021-04484-y

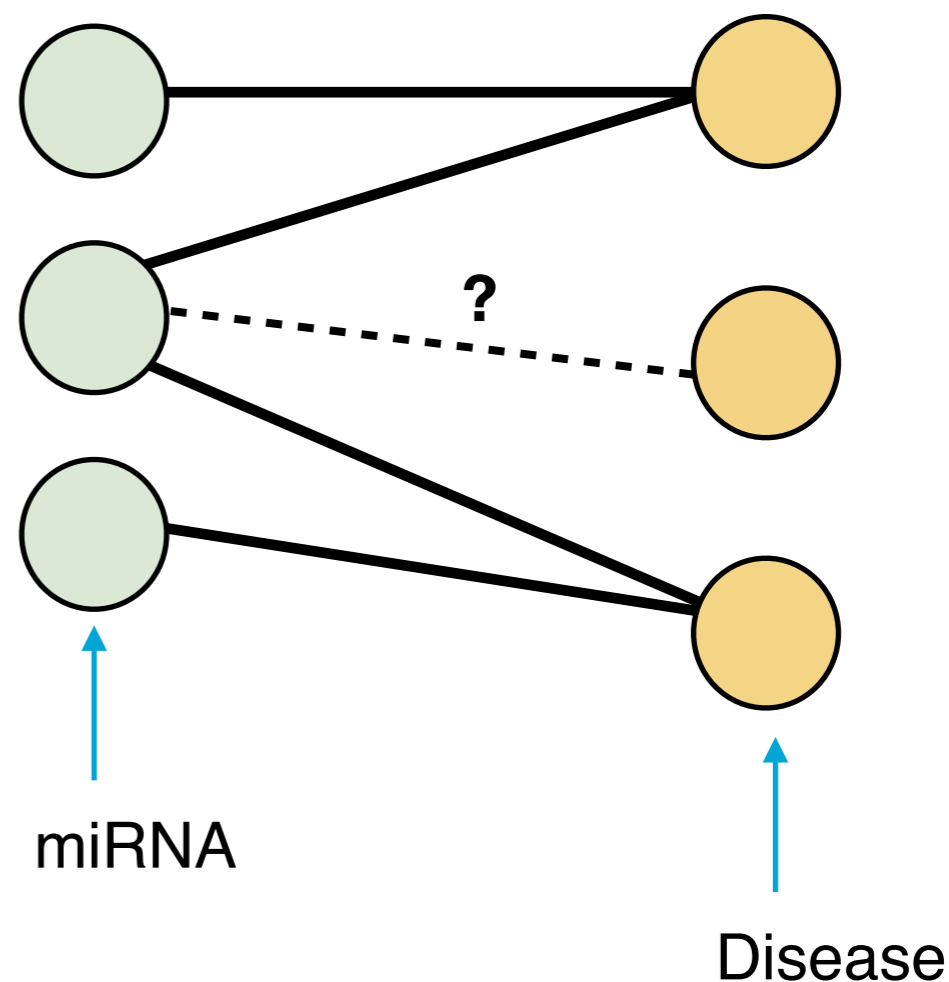# Predicting miRNA-disease associations

Data bias

20% of the diseases account for 80% of associations
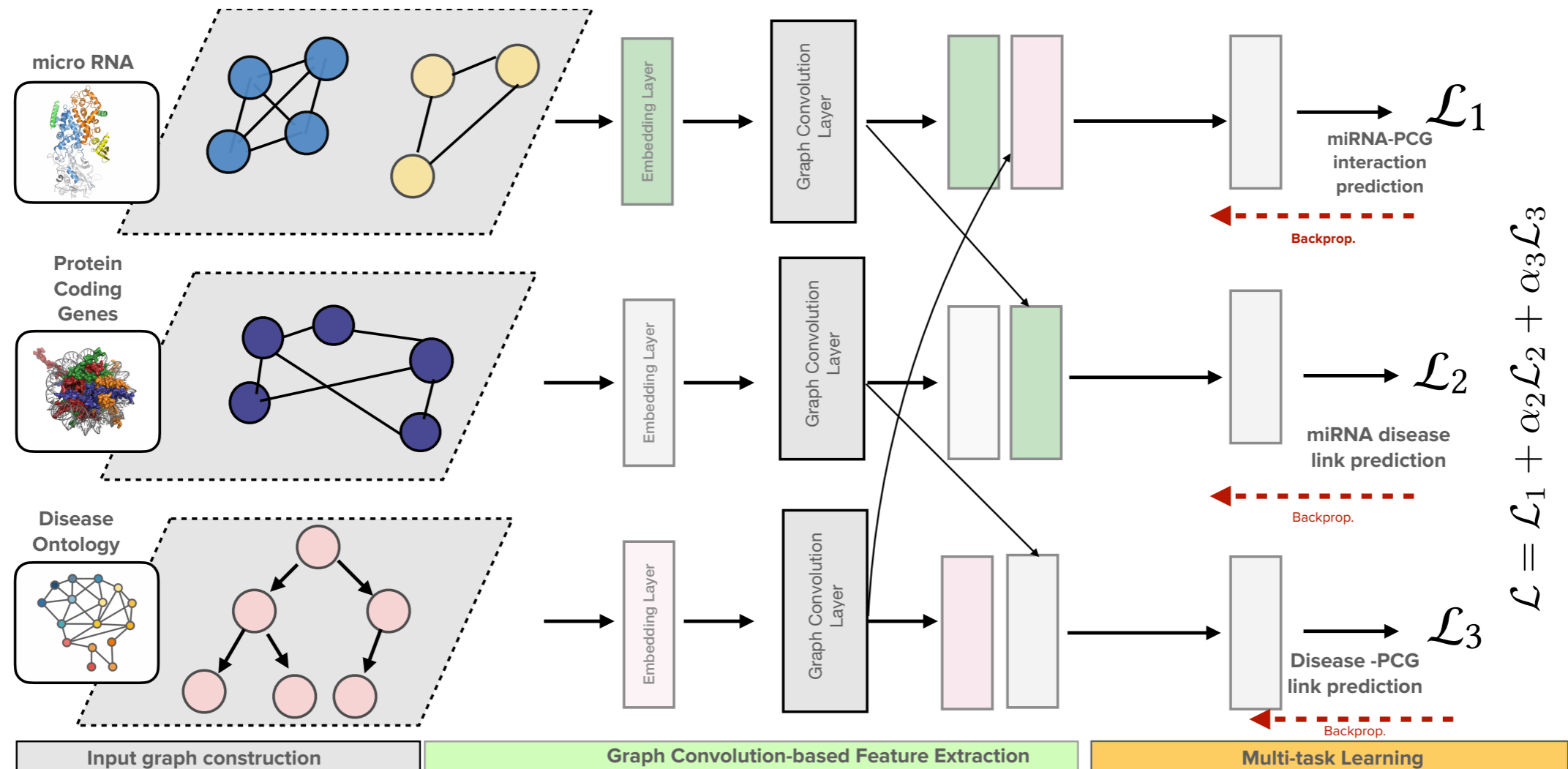
Data scarcity

Sparse bipartite graph with small number of nodes

High number of false positives in training data

?

miRNA

Disease

Overall strategy: Learn jointly from miRNA family, miRNA-gene, disease-gene interactions and disease ontology information
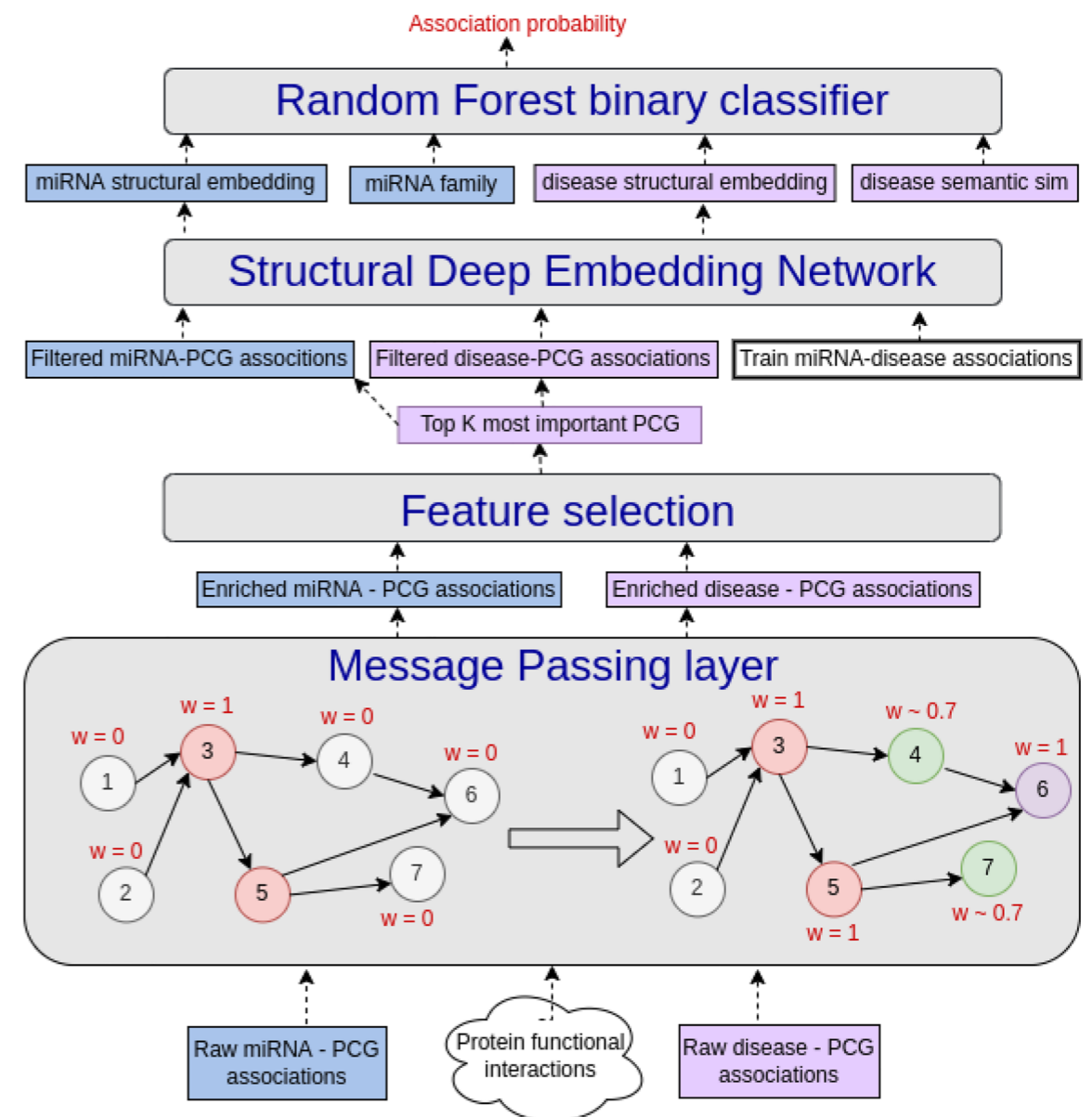
# Join learning framework

# How to filter training data

- Message passing to enrich gene associations

- Feature selection to select important genes

- Shallow network embeddings over Heterogeneous association network
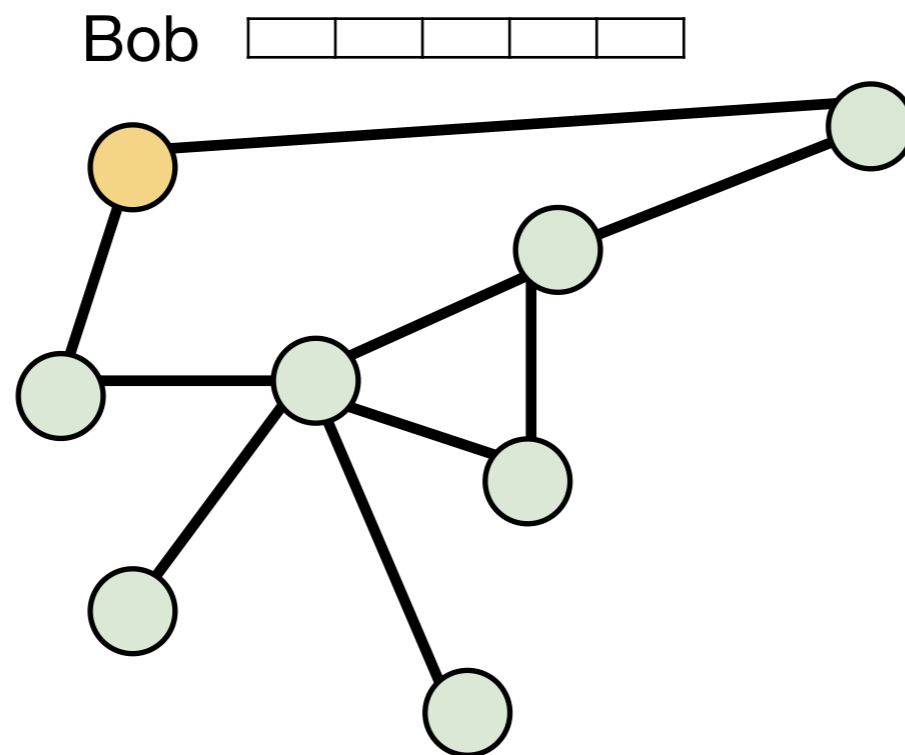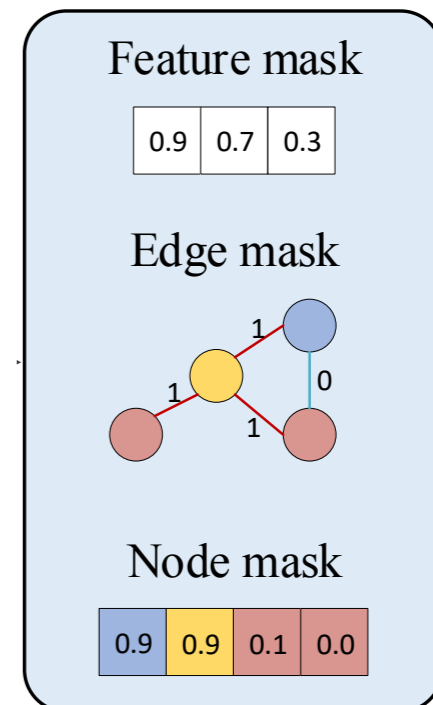


https://www.nature.com/articles/s41598-022-20529-5

# Challenges of transparency and privacy

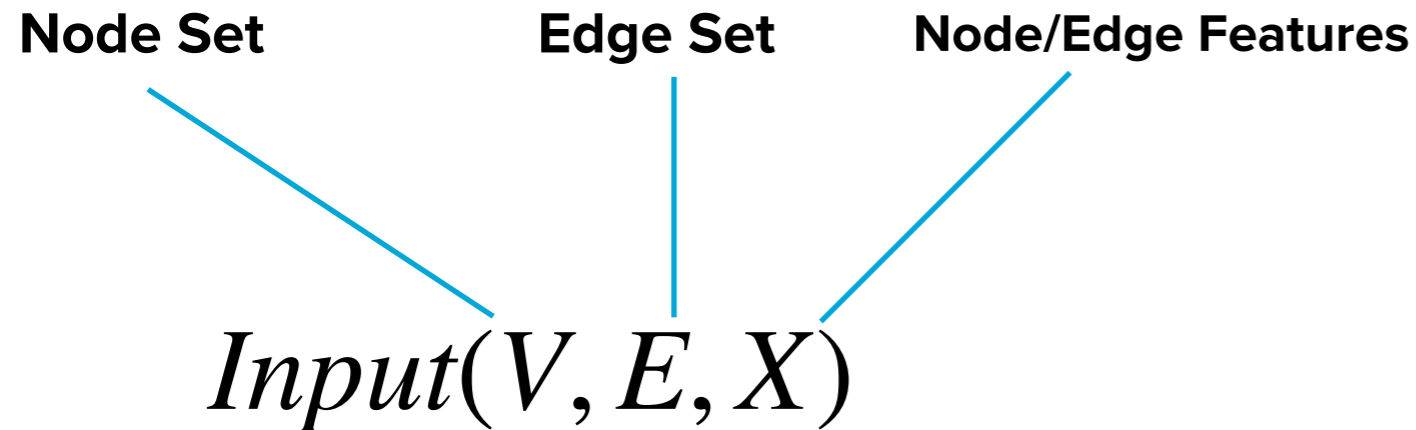# Transparency

Why was Bob's loan denied?
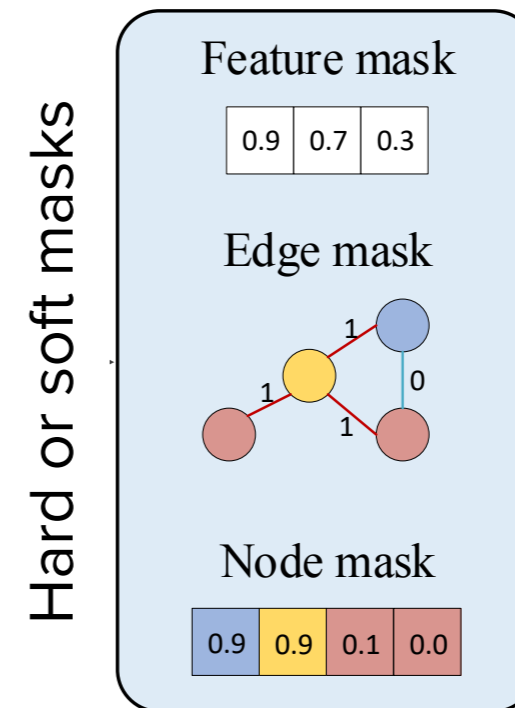
**Explanation**



Decision has to be explained not only in terms of features but also graph structure. General explainability methods cannot be trivially applied for graphs.

# Post-hoc explanations

**Node Set**     **Edge Set**     **Node/Edge Features**

$$Input(V, E, X)^A$$

$X$

**Explanation types**



Hard or soft masks

Feature mask

| 0.9 | 0.7 | 0.3 |

Edge mask

Node mask

| 0.9 | 0.9 | 0.1 | 0.0 |

Examples: GNNExplainer, Zorro, PGExplainer

**Explanation types:**

Feature explanations in terms of most relevant features $X' \subset X$

Structure explanations in terms of most relevant nodes $(V' \subset V)$ or edges $(E' \subset E)$

We are interested in finding both feature and structure explanations which effectively capture interplay of structure and features in model's decision making.
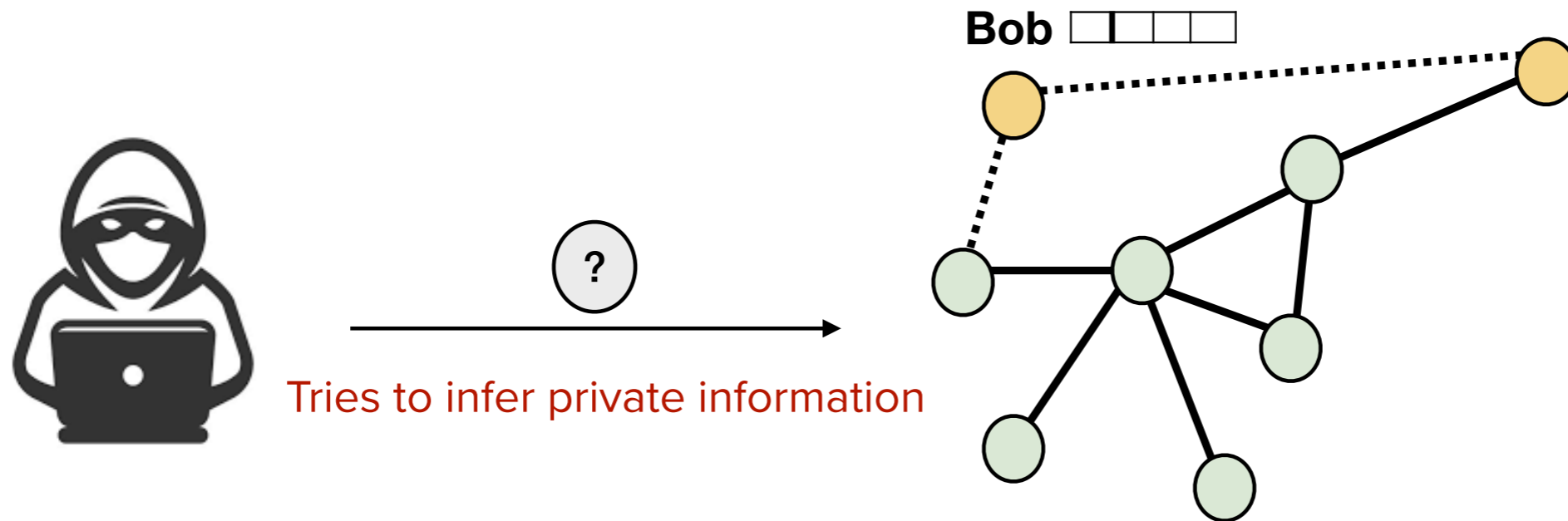
# Privacy

**Graphs can contain sensitive information**

- User's sensitive attributes
- Sensitive relations

**GNNs encode relation information within the model, could memorise such information**

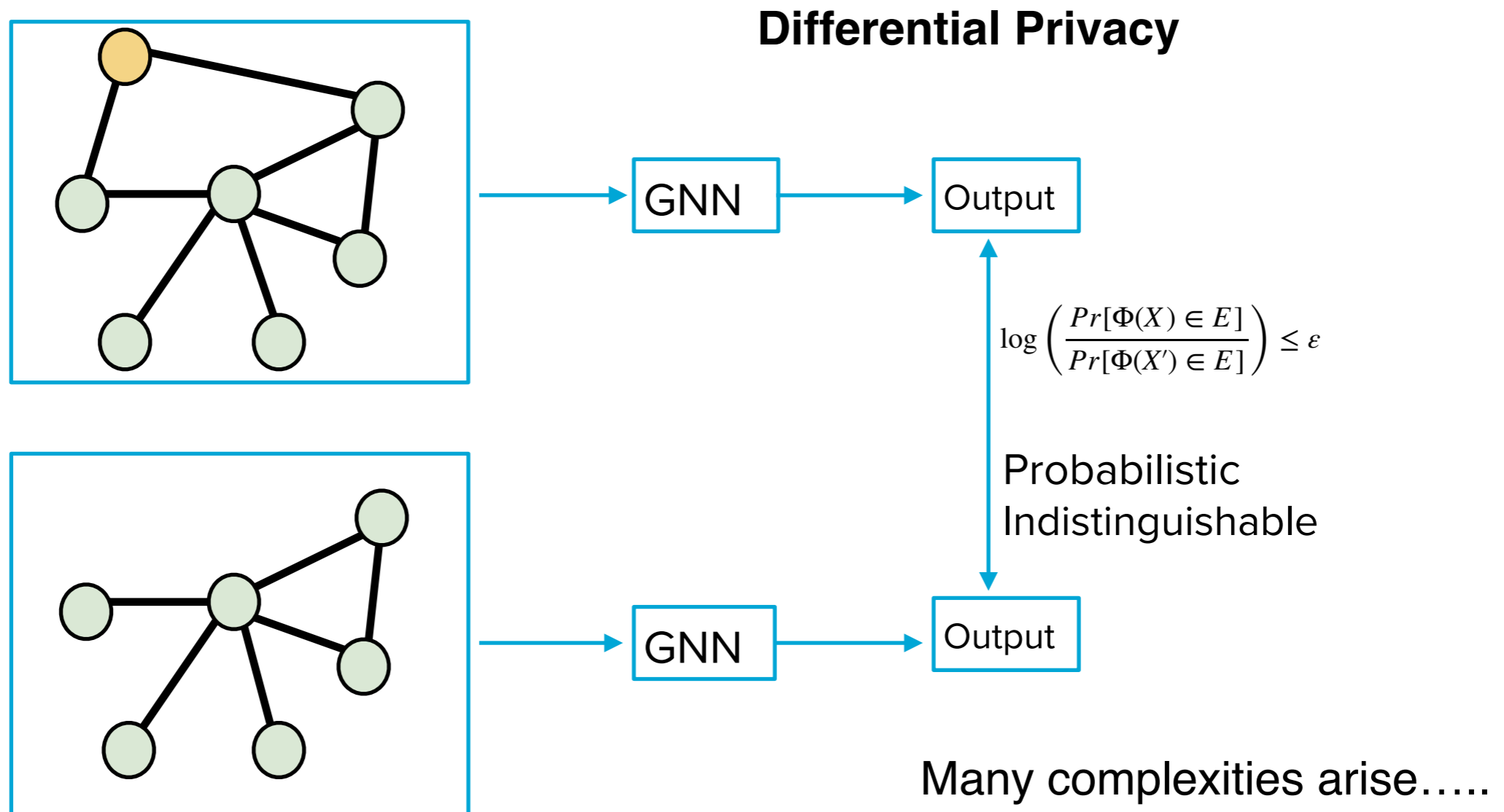- Your identity could be revealed because of your neighbour

# Privacy



**Bob**

? 

Tries to infer private information

**Node Membership Inference :** Is Bob a part of training data? [Olatunji et al., '21] [Duddu et al., '20]

**Relation reconstruction :** Who are friends of Bob? [He et al., '21] [Zhang et al., '20]

**Attribute Inference :** Does Bob smoke?

# Building Private GNN Models



**Differential Privacy**

$$\log\left(\frac{Pr[\Phi(X) \in E]}{Pr[\Phi(X') \in E]}\right) \leq \varepsilon$$

Probabilistic
Indistinguishable

Many complexities arise…..

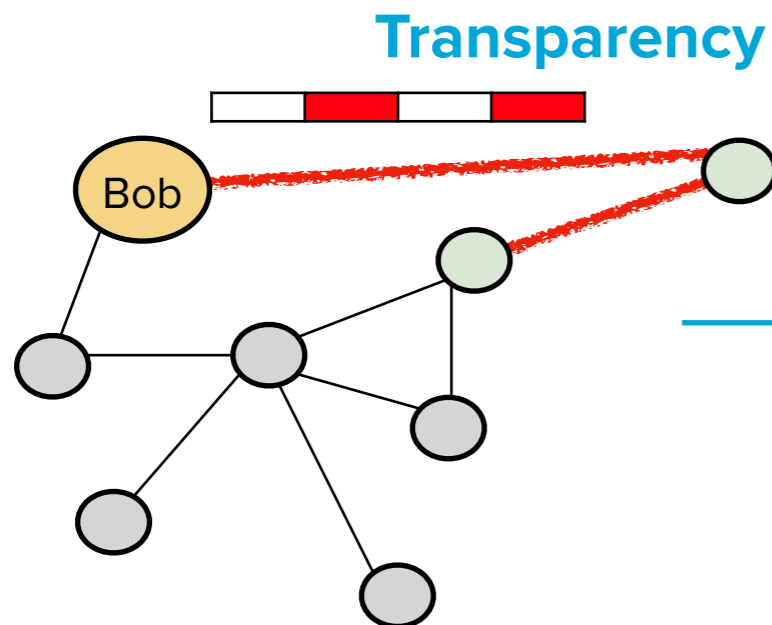# Building Privacy Preserving models for graphs

**A direct application of techniques like DP-SGD is not possible due to**

**-** Unbounded sensitivity (think of the effect of leaving out or adding one node in a graph)

**-** Violation of i.i.d. assumption

**-** Need for inference privacy (as training data might be used during inference)

PrivGNN (Olatunji, Funke, Khosla, 2021), GAP (Sajadmanesh, Shamsabadi, A, Bellet,  et al. 2022)
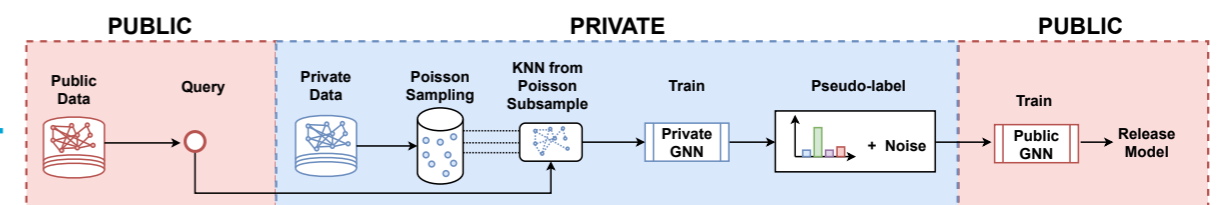
**TU**Delft

# Transparency - Privacy Tradeoffs

But we want our models to be **transparent** and
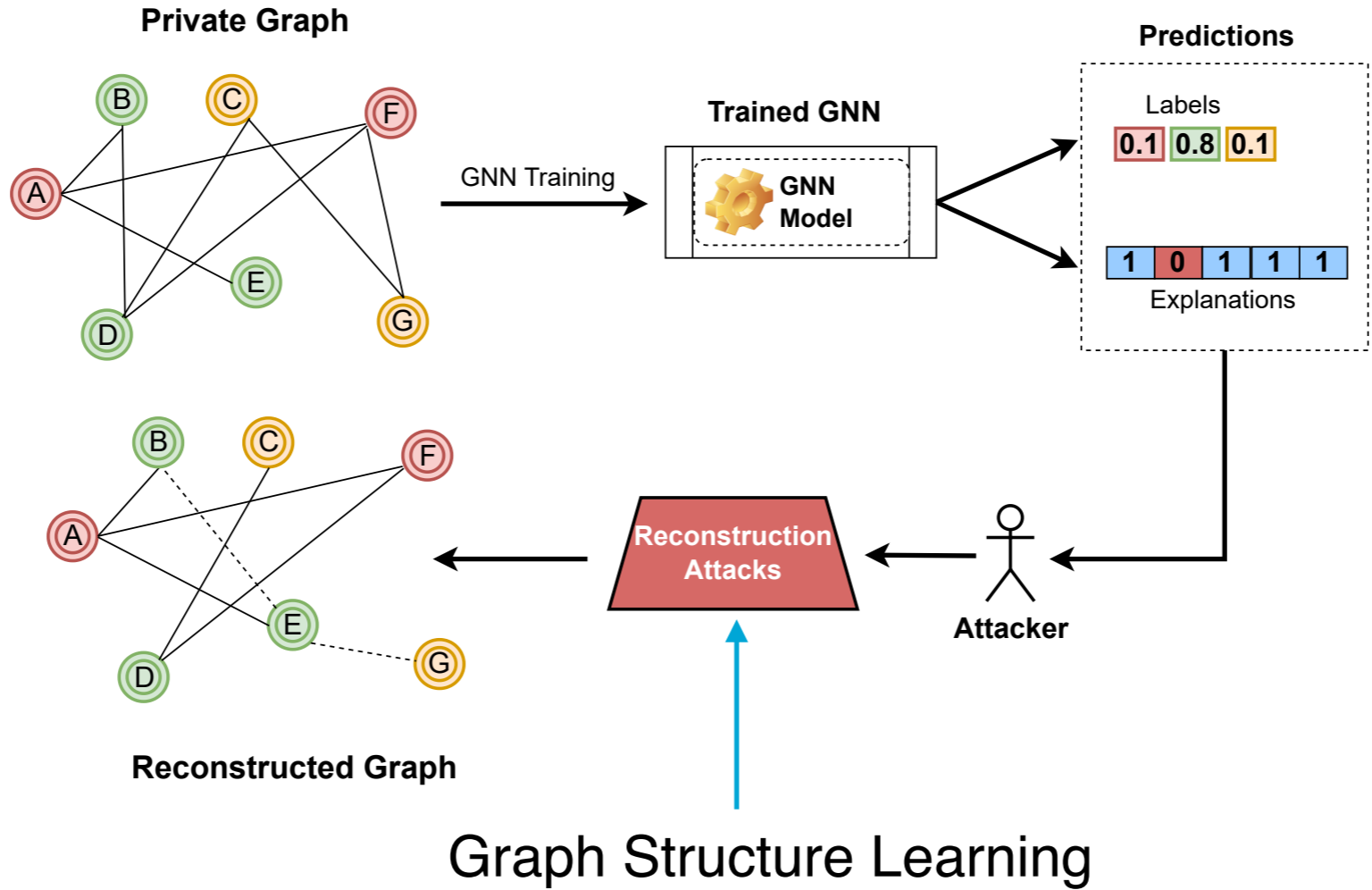**private** simultaneously

**Transparency**



**Privacy**

**How much can the privacy be hurt?
How to preserve privacy?**

**Can we explain private
models? How?**



**Complex privacy preserving mechanism**

# Reconstructing graphs from feature explanations



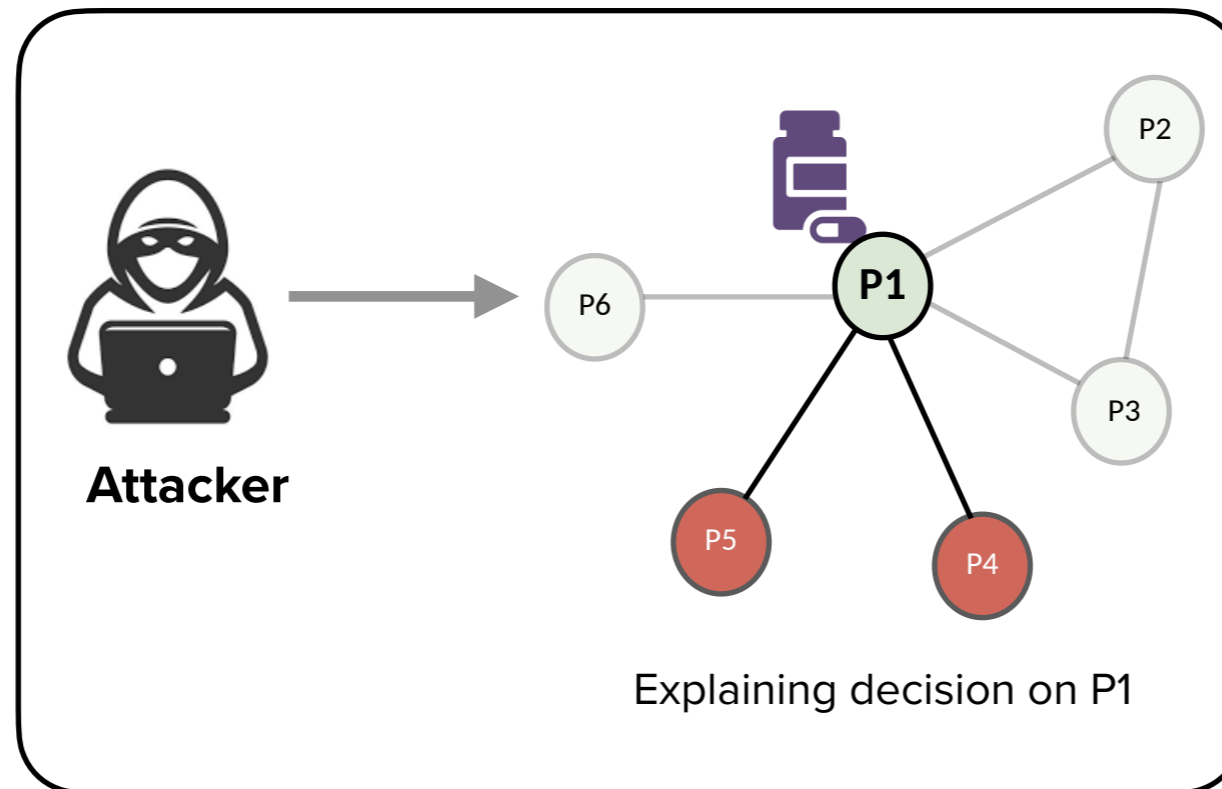Private Graph Extraction via Feature Explanations [Olatunji et al. 2022]

https://arxiv.org/abs/2206.14724

# Some interesting findings

**-** Training graph could be reconstructed using alone the feature explanations and the labels

**-** Certain explanations leak more information than others

**-** Gradient based explanations incur high privacy loss while showing low **utility** (quantified by high faithfulness and sparsity)
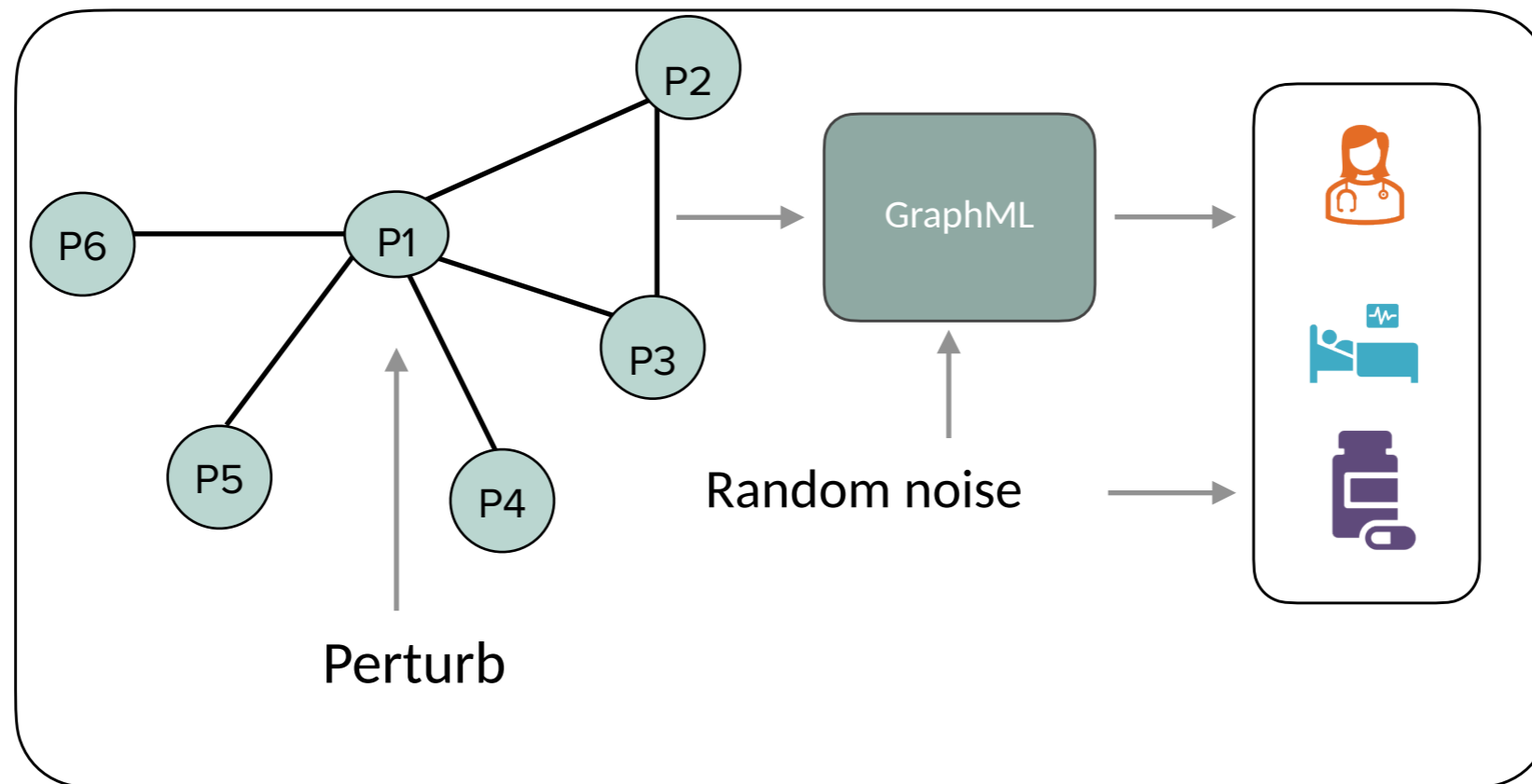
# Challenges

Structure explanations can directly reveal information about neighbours



Explaining decision on P1

Explanations of neighbouring datapoints would be correlated

# Challenges

Private learning over graphs is more complex than that in standard ML



How to define explanation for a private model?

# Research Directions and Open Questions

**Quantification of privacy leakage in presence of different explanation types**

- How can we measure information leakage due to different explanation types?
- Risk-utility assessment of different explainers/explanations
- Can we release explanations privately while still maintaining their utility?

**Explaining the decisions of privacy-preserving models**

- What should be the properties of an explanation for a privacy-preserving model?
    - Such properties might need to be defined based on the private learning strategy
- How to release such explanations in a private manner?

**Joint optimization of privacy and transparency**

- How can we optimise for the combined requirements of privacy and transparency in GraphML?

**Privacy and Transparency in Graph Machine Learning: A Unified Perspective**, M.Khosla. In AIMLAI@CIKM'22

https://arxiv.org/abs/2207.10896