

# Link Stealing Attacks Against Inductive Graph Neural Networks

*Yixin Wu,<sup>1</sup> Xinlei He,<sup>2</sup> Pascal Berrang,<sup>3</sup> Mathias Humbert,<sup>4</sup>  
Michael Backes,<sup>1</sup> Neil Zhenqiang Gong,<sup>5</sup> Yang Zhang<sup>1</sup>*

<sup>1</sup> CISPA Helmholtz Center for Information Security

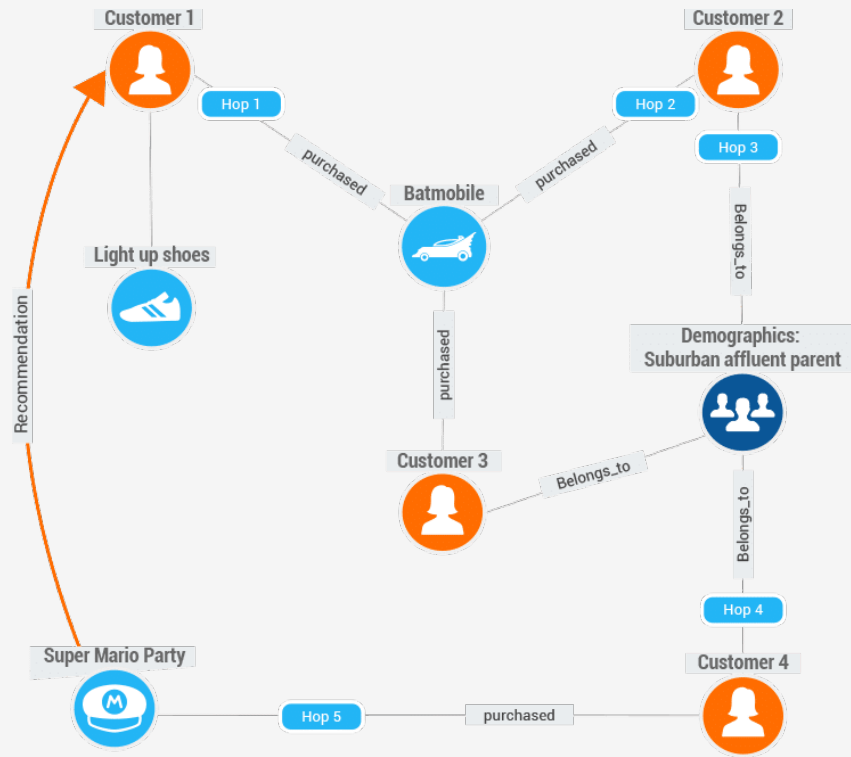
<sup>2</sup> Hong Kong University of Science and Technology

<sup>3</sup> University of Birmingham <sup>4</sup> University of Lausanne <sup>5</sup> Duke University

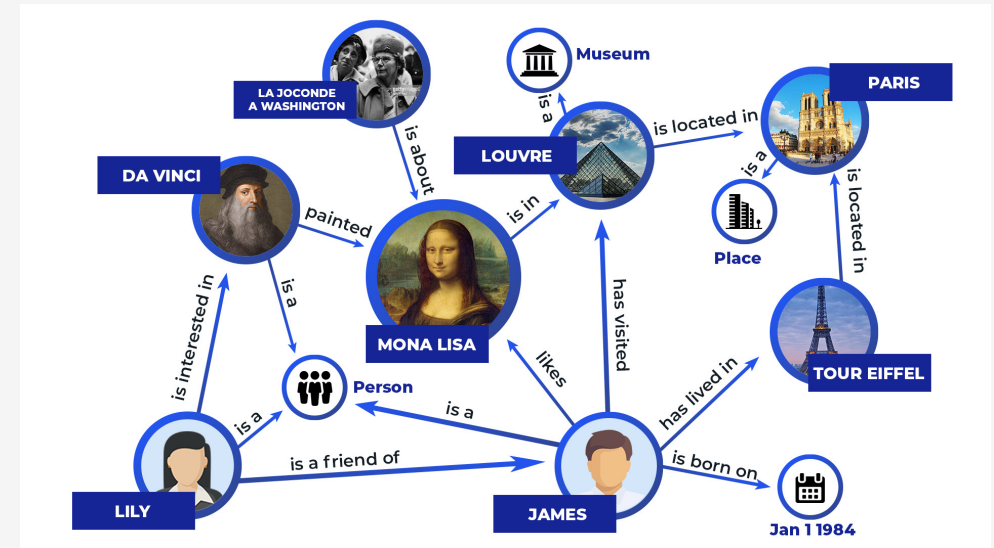




# Prevalence of Graph-Structured Data



## Recommendation System



## Knowledge Graph



## Social Network

Image source: <https://www.tigergraph.com/solutions/recommendation-engine/>;

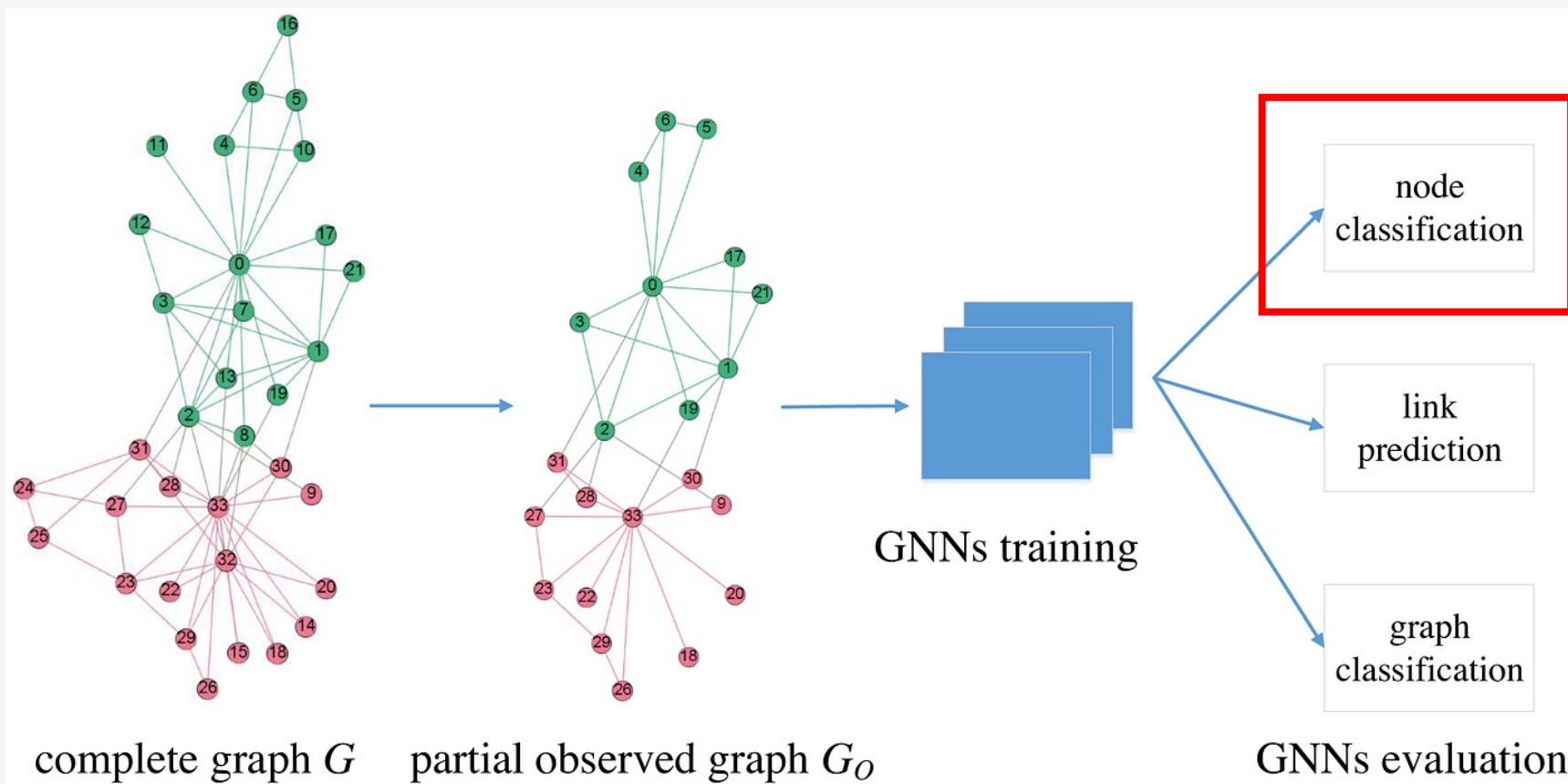
2 <https://yashueth.wordpress.com/2019/10/08/introduction-question-answering-knowledge-graphs-kgqa/>

<https://medium.com/analytics-vidhya/social-network-analytics-f082f4e21b16>



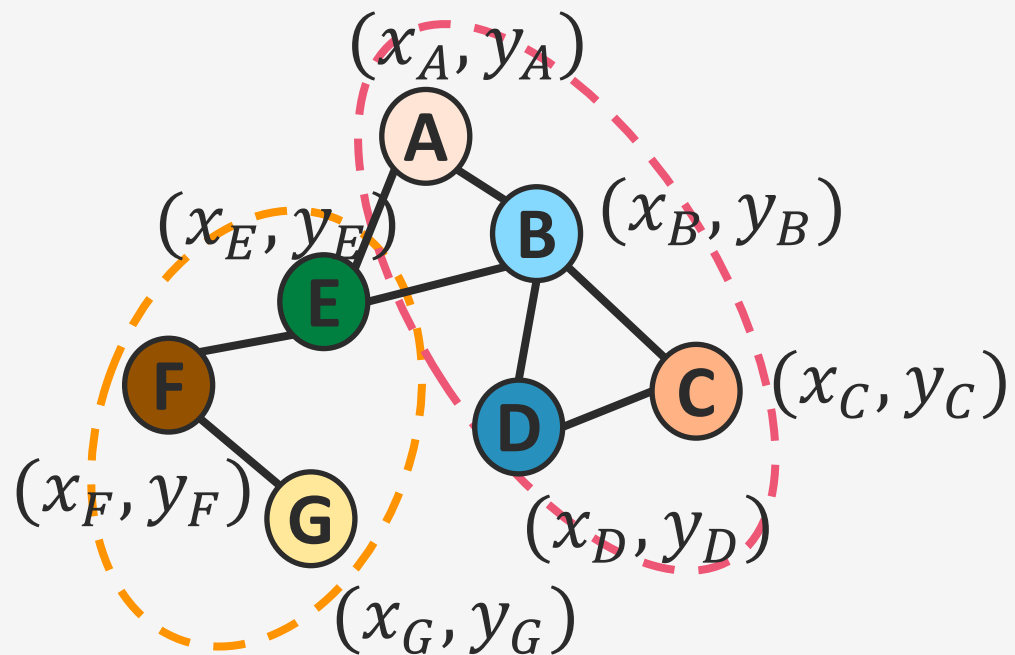
# Graph Neural Networks (GNNs)

- Leverage neighbor information among nodes to learn embeddings to perform downstream tasks





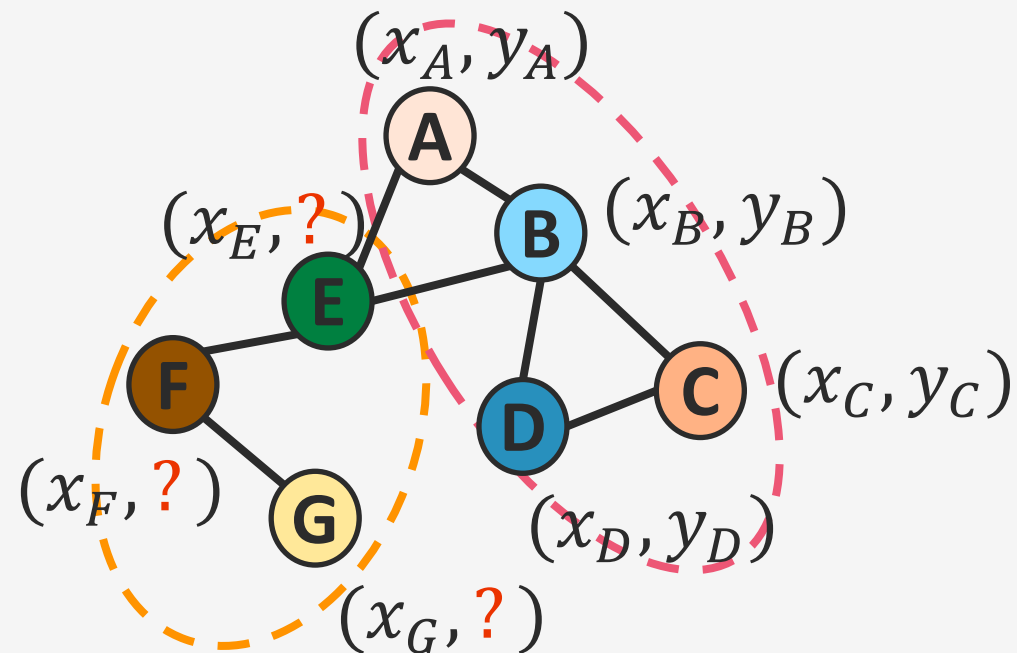
# Transductive Setting



## Train/Test Split

During training, **the entire graph** including node attributes and edges can be observed

## Training

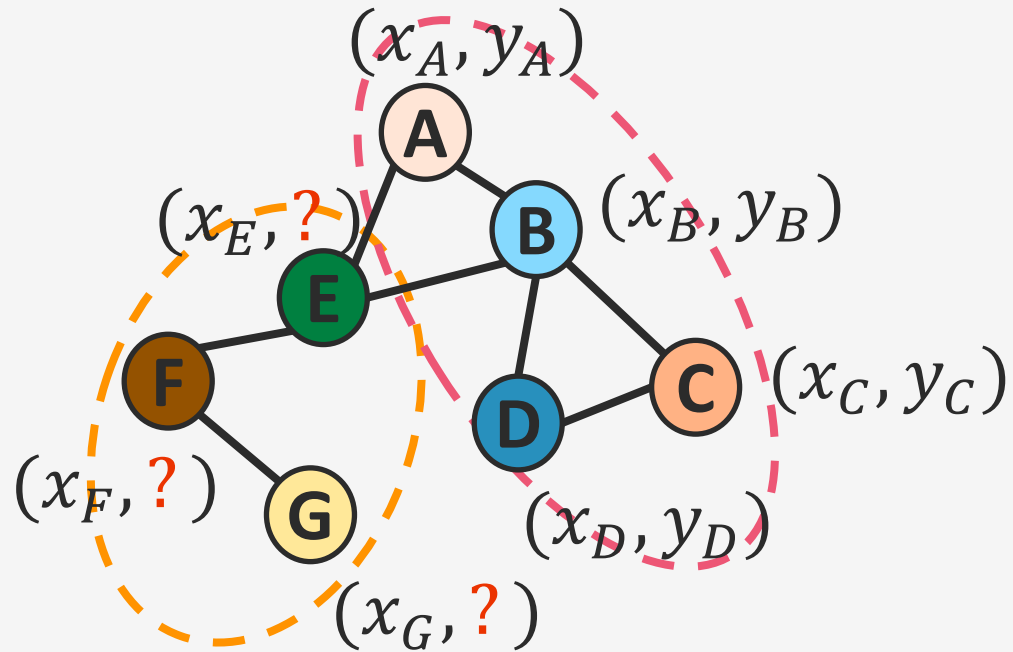


Transductive  
GNN



# Transductive Setting

## Graph Observed During Training



## Testing

G



Transductive  
GNN

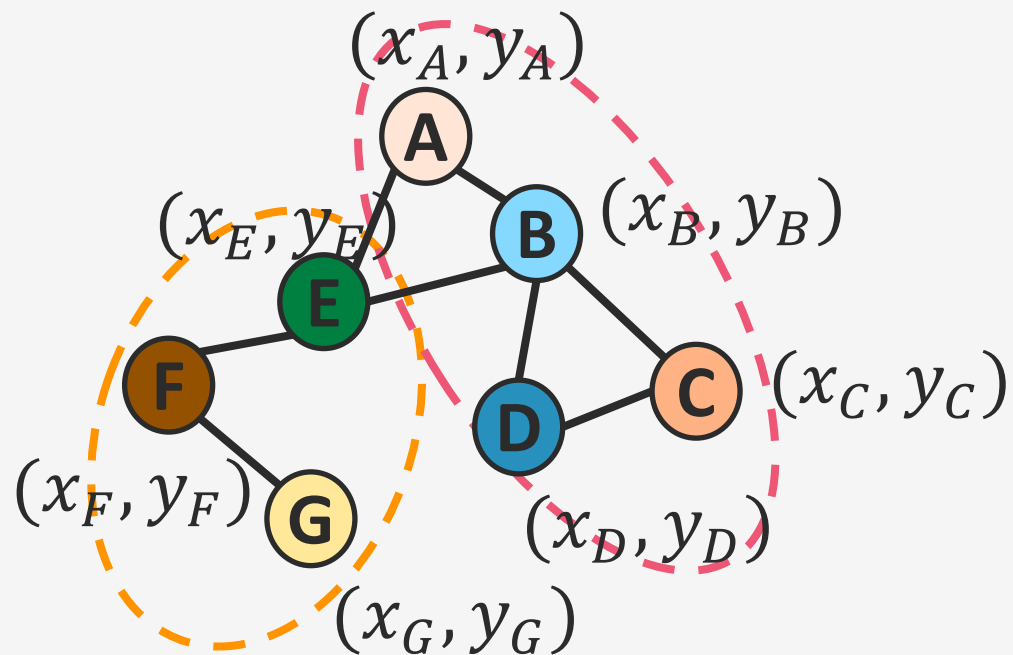


$y_G$

During inference, users feed the **identifiers of unlabeled seen nodes** into GNN to obtain prediction results



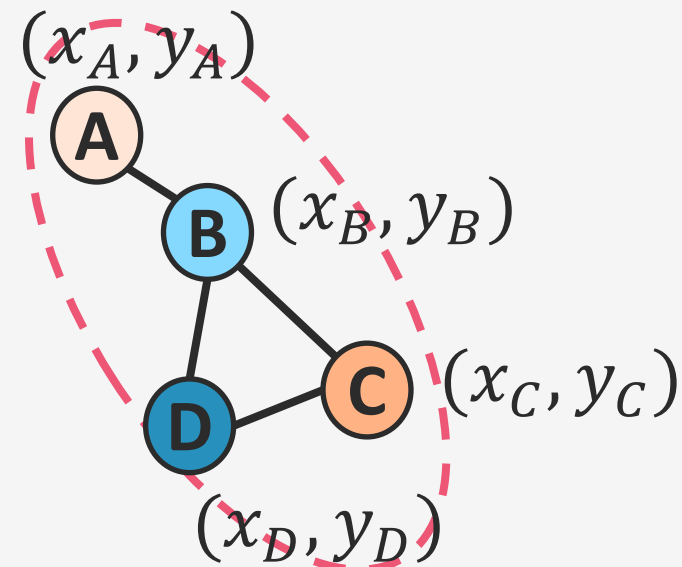
# Inductive Setting



## Train/Test Split

During training, **only the training graph** can be observed

## Training



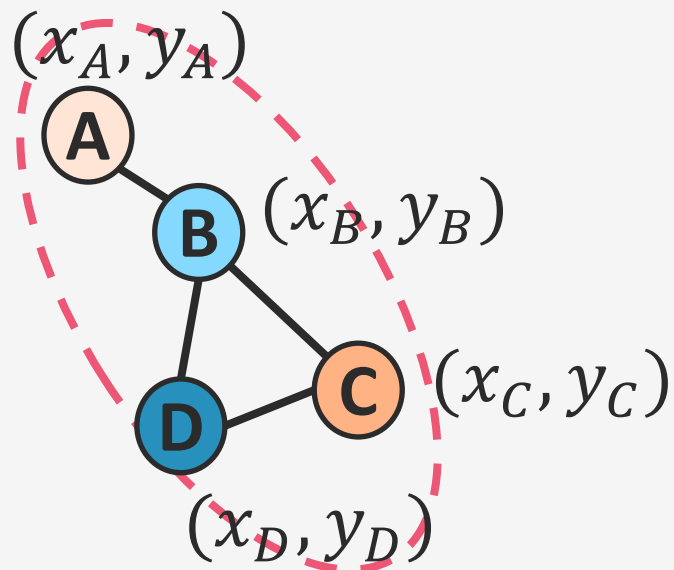
Inductive  
GNN



# Inductive Setting

## Testing

### Graph Observed During Training

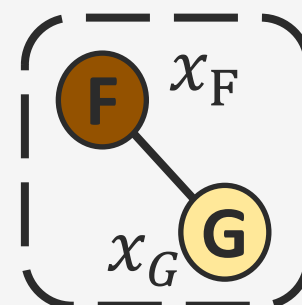


### 0-hop



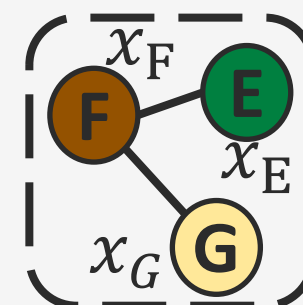
or

### 1-hop



or

### 2-hop



Inductive  
GNN



$y_G$

During inference, users **construct subgraphs** to obtain prediction results of **unseen nodes**



## Difference Between Two Settings

**Transductive  
GNN**

“Memorize” the training graph

**Inductive  
GNN**

Learn a generalizable embedding function

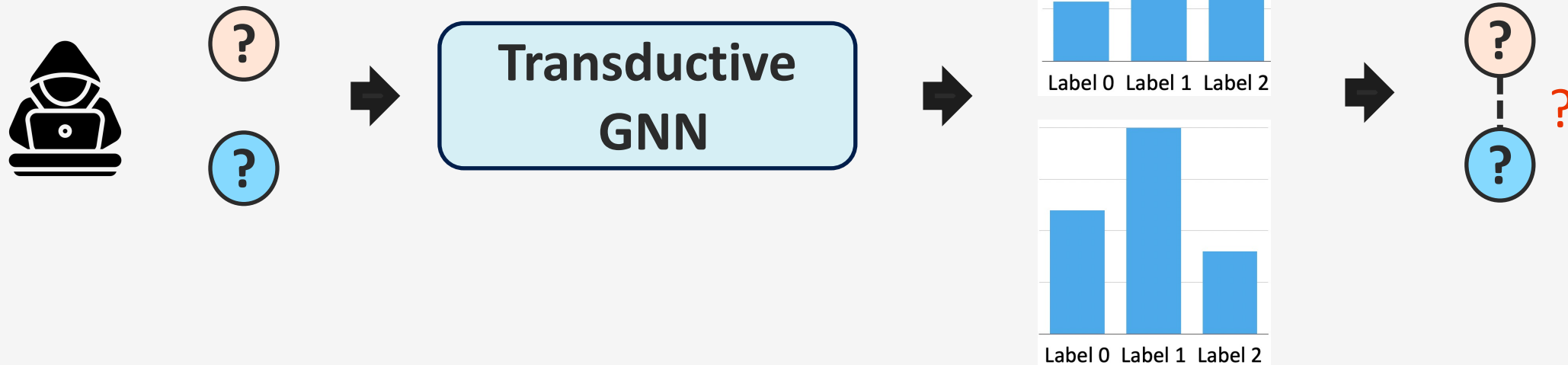
Inductive GNNs are more generalizable and flexible for dynamic real-world applications, e.g., social network and recommendation system





# Link Stealing Attacks on Transductive GNNs

- Previous work<sup>[1]</sup> demonstrates that the transductive GNNs are vulnerable to link stealing attacks



Given two nodes used to train a black-box GNN, can we predict whether they are linked?



# Challenges From The Differences

**Transductive  
GNN**

“Memorize” the training graph

**Inductive  
GNN**

Learn a generalizable embedding function

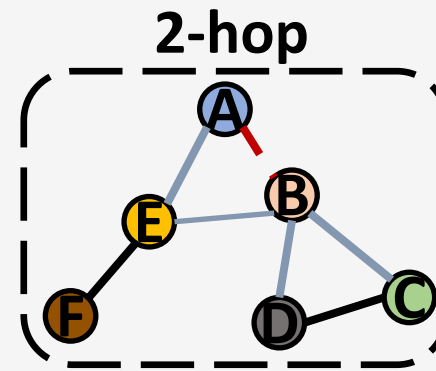
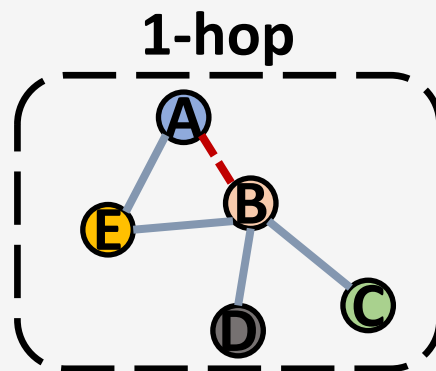
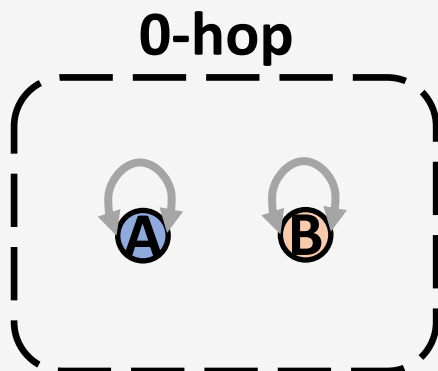
Inductive GNNs might not learn specific information of training graph as transductive GNNs



# Challenges From The Differences

**Inductive  
GNN**

Query with own constructed subgraphs



The adversary relies on limited and incomplete neighbor information, as the information of the link they intend to infer is missing



# Are inductive GNNs vulnerable?

**Transductive  
GNN**

Query with the **identifiers of unlabeled nodes** and obtain fixed node embeddings learned during training

**Inductive  
GNN**

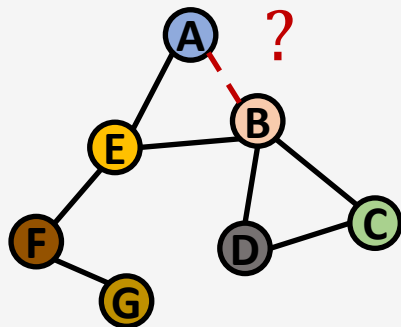
Query with own constructed subgraphs and obtain node embeddings based on the subgraphs

Given the above two challenges, are inductive GNNs vulnerable to link stealing attacks?

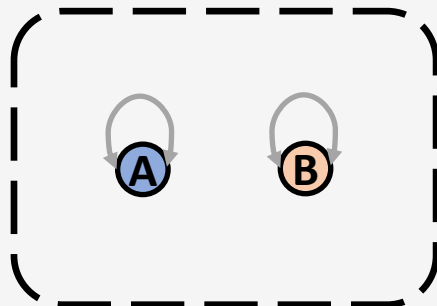


# Link Stealing Attacks on Inductive GNNs

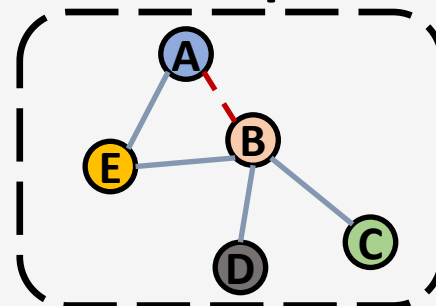
Adversary



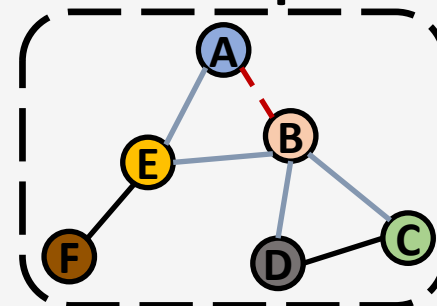
0-hop



1-hop

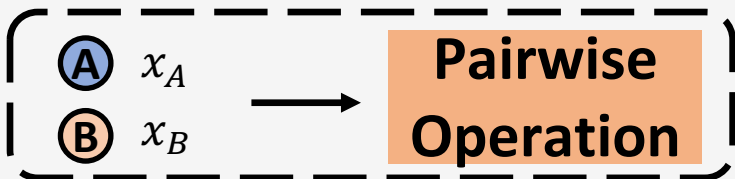


2-hop

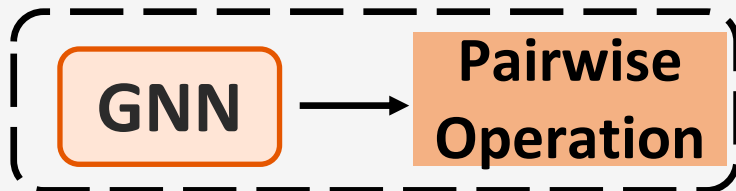


Adversary can have either these three types of features

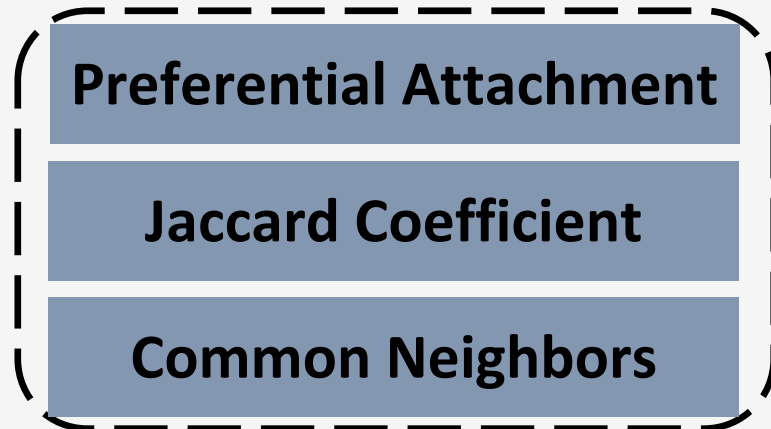
Node Features



Posterior Features



Graph Features





# Evaluation Results

Dataset	Posterior-Only Attack			Combined Attack						
	A0	A1	A2	A3	A4	A5	A6	A7	A8	A9
Cora	0.859	0.849	0.849	0.876	0.876	0.875	0.882	0.884	0.908	<b>0.909</b>
Pubmed	0.768	0.806	0.809	0.889	0.895	0.897	0.881	0.882	<b>0.939</b>	<b>0.939</b>
DBLP	0.781	0.821	0.822	0.834	0.873	0.872	0.879	0.903	0.924	<b>0.929</b>
Photo	0.877	0.898	0.898	0.892	0.916	0.915	0.967	<b>0.968</b>	0.946	0.946
CS	0.817	0.838	0.845	0.869	0.890	0.893	0.955	<b>0.956</b>	0.941	0.940
LastFM	0.850	0.869	0.867	0.883	0.909	0.911	0.919	0.921	0.929	<b>0.930</b>

- The proposed attacks with no (A0; 0-hop only) or limited (A1; 1-hop query) neighbor information can achieve good performance
- More information achieves better performance



# Conclusion

- We propose in total 10 link stealing attacks against inductive GNNs
- No neighbor information (0-hop query) still enables well-performing link stealing attacks
- More information achieves better attack performance
- High robustness of the proposed attacks; better performance than traditional link prediction (baseline), showing inductive GNNs indeed leak privacy information



# Thanks!



**Yixin Wu**

**CISPA Helmholtz Center for Information Security**

**@yxoh28**

**<https://yxoh.github.io/>**