

Graph Data Mining and Learning: Modern applications, challenges, and paradigms

Naheed Anjum Arafat

Post-doc

Nanyang Technological University, Singapore (2021-2024)

PhD

National University of Singapore

January 22, 2025

Outline

1 Introduction

2 My Research

- Learning on Graphs: What, When, and Why

3 Key Research Themes

- Modern Application of Graph Learning
- Modern Challenges to Graph Learning
- Graph Data Mining and Learning: Modern paradigms

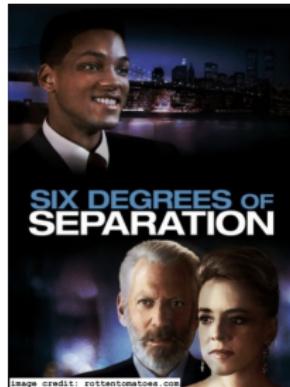
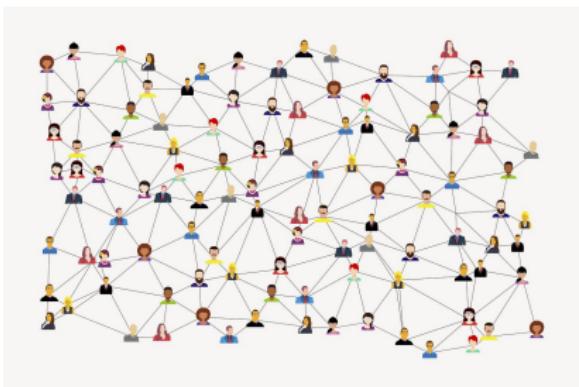
4 Future Plans

5 Interdisciplinary Works and Collaborations

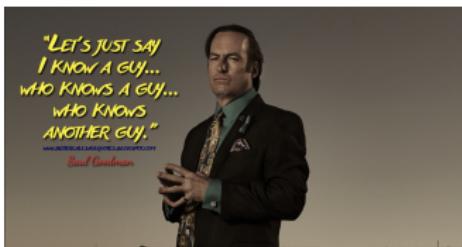
6 Mentoring

7 Q&A

Graphs before..



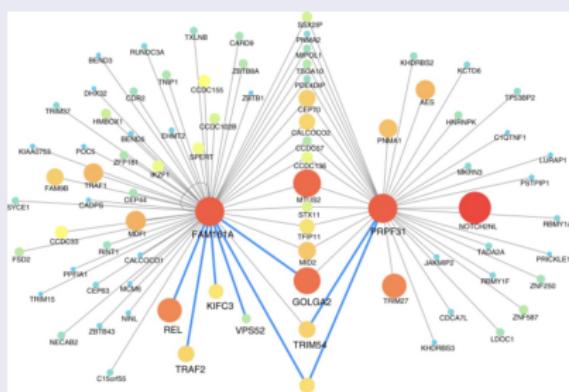
Graphs before..



Graphs now..

Well, they are everywhere

Biology



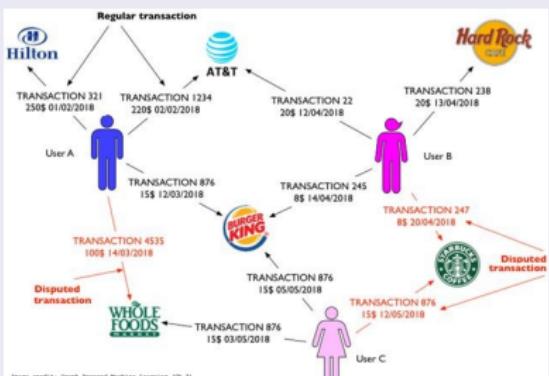
Protein-protein interactions

Image credit: Kovács, István A., et al. "Network-based prediction of protein interactions." *Nature comm.*, 2019

Graphs now..

Well, they are everywhere

Finance

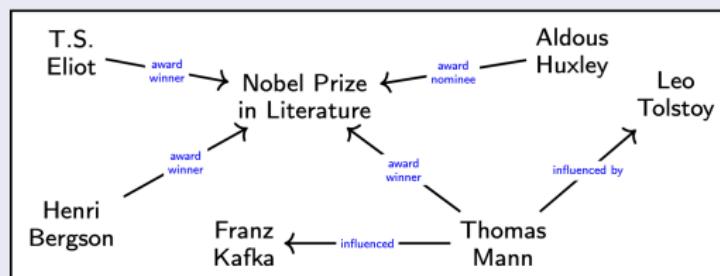


Financial transactions

Graphs now..

Well, they are everywhere

Human centered computing



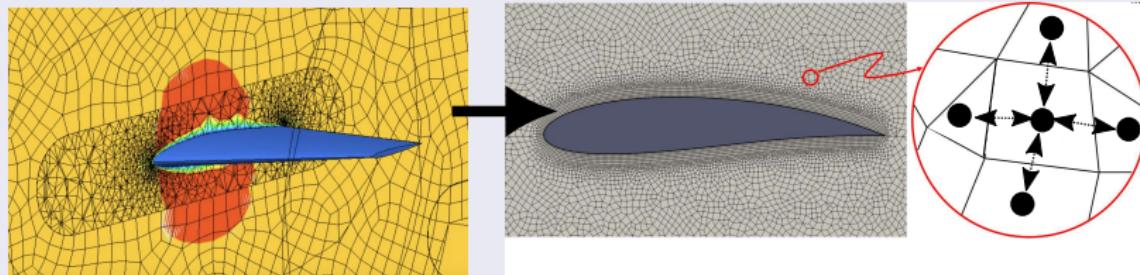
Knowledge graphs for explainable AI

Modern paradigm: Human-understandable explanation to black-box AI models. For instance, Recommending *War and Piece* because you have read *The Magic Mountain* by Thomas Mann who was influenced by works of Tolstoy.

Graphs now..

Well, they are everywhere

Physics



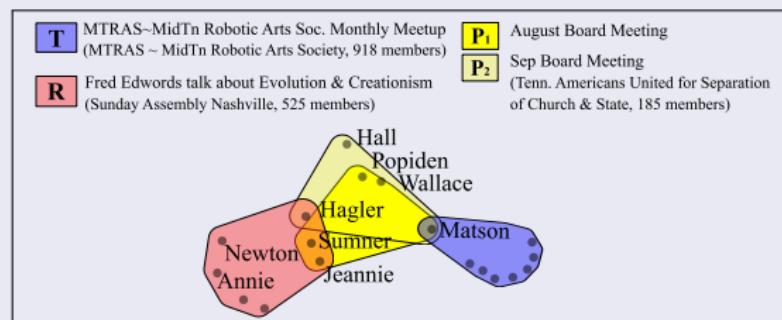
Fluid simulation at cross-section of an airplane wing (aka airfoil) and a close-up view

Modern paradigm: Modeling fluid flow as message passing on graphs.

Graphs now..

Well, they are everywhere

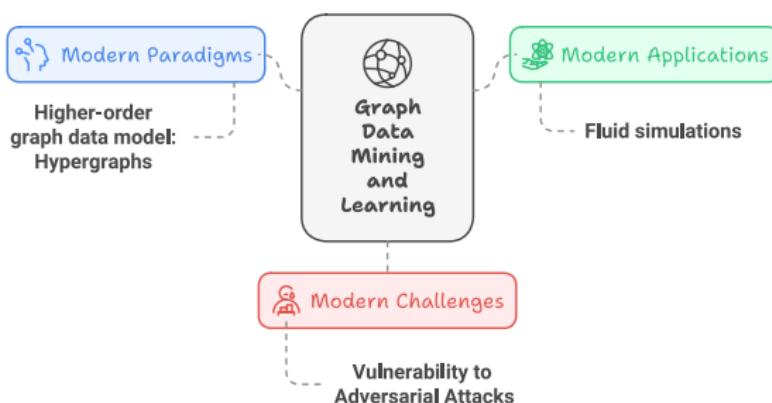
Public Health and Epidemiology



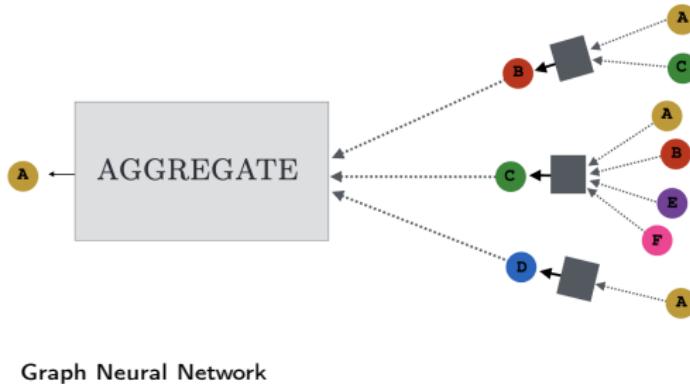
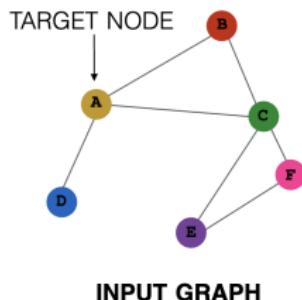
Modern paradigm: Modeling higher-order interactions as Hypergraphs.

My Research: Overview

My research aims to advance our capabilities in extracting insights from graphs and learning on graphs.



My Research: Learning on Graph

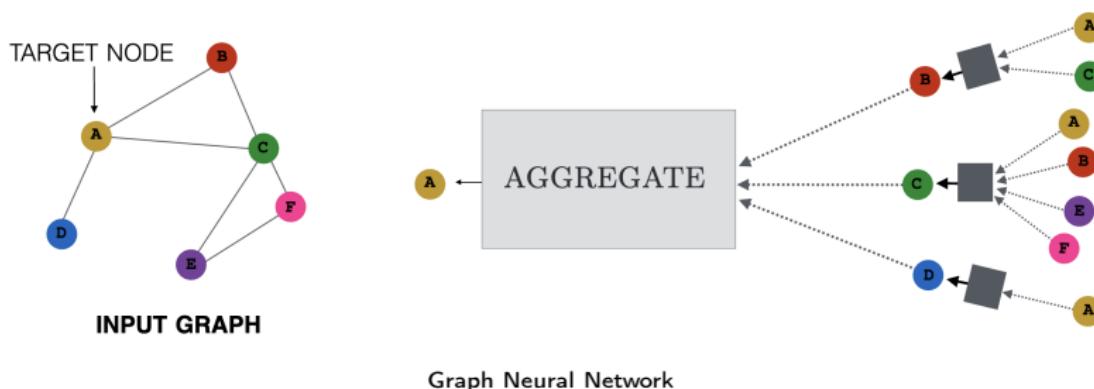


When GNNs (applications):

- Node classification (identifying fraudulent transactions)
- Link prediction (recommending movies, videos, posts, products)
- Graph classification (identifying toxic chemicals)

Basically, any learning task involving entities connected by relations.

My Research: Learning on Graph

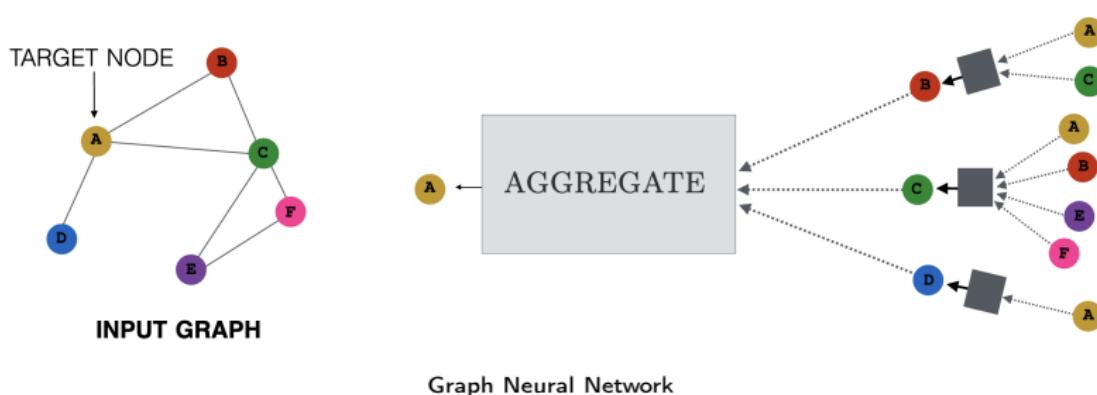


When GNNs (applications):

- Node classification (identifying fraudulent transactions)
- Link prediction (recommending movies, videos, posts, products)
- Graph classification (identifying toxic chemicals)

Basically, any learning task involving entities connected by relations.

My Research: Learning on Graph

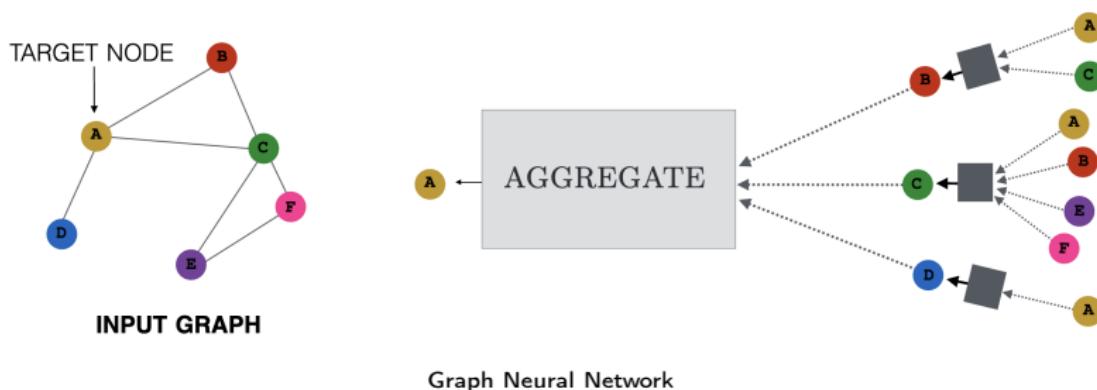


When GNNs (applications):

- Node classification (identifying fraudulent transactions)
- Link prediction (recommending movies, videos, posts, products)
- Graph classification (identifying toxic chemicals)

Basically, any learning task involving entities connected by relations.

My Research: Learning on Graph



When GNNs (applications):

- Node classification (identifying fraudulent transactions)
- Link prediction (recommending movies, videos, posts, products)
- Graph classification (identifying toxic chemicals)

Basically, any learning task involving entities connected by relations.

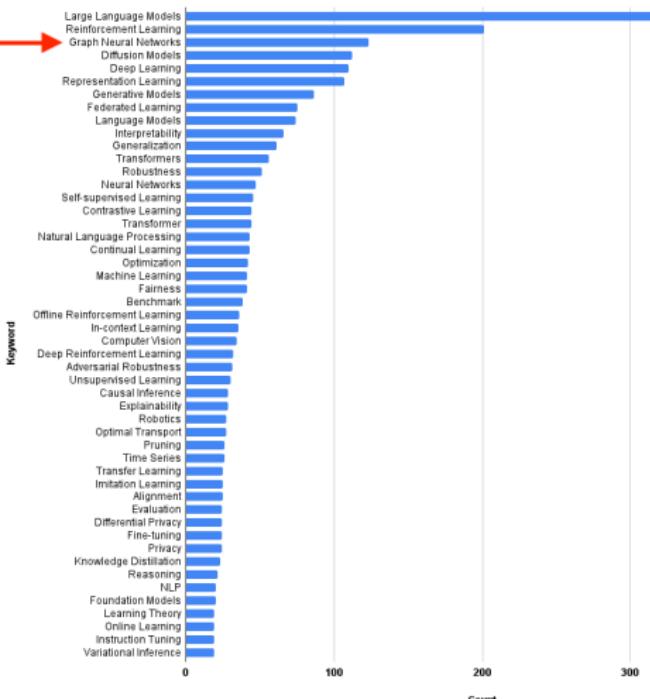
My Research: Why GNNs?

GNNs in Academia: Top 50 keywords in submitted research papers at ICLR 2024.

GNNs in Industry:

- Uber (food/restaurant recommendations)
- Alibaba (product recommendation)
- Snapchat (content recommendations)
- Pinterest (Pins recommendation)
- Google (Google Maps and Deepmind)
- Kumo.ai (fraud detection)
- Ant financial Services (fraud detection)
- Meta (preventing spread of misinformation, fake account detection)

50 Most Frequent Keywords



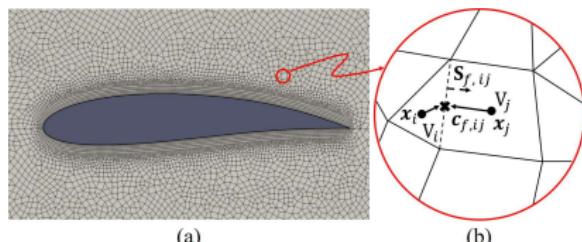
Fluid Simulation: a Modern Application of GNNs

Finite Volume Features, Global Geometry Representations, and Residual Training for Deep Learning-based CFD Simulation, ICML 2024 (spotlight)

- **Problem:** Predict velocity, pressure at every point in the domain using AI.

- **Motivation:**

- Numerical simulations of fluids are computationally expensive often requiring days.
- Recently, data-driven GNNs have been deployed as an efficient alternative. (e.g. MeshGraphNet from Deepmind)
- Can we improve the fidelity of data-driven GNNs?



(a) CFD mesh with an airfoil body surrounded by different sizes of cells. (b) Illustration of cell characteristics, namely cell centroids, x_i and x_j , face centroid, $c_{f,ij}$, face area normal vector, $S_{f,ij}$, and cell volumes, V_i and V_j .

Contributions

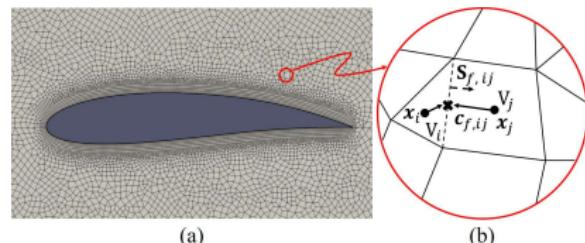
- 1 Geometric features to inform GNN nodes of long-range interactions.
- 2 Finite-volume features in the graph convolution layer.
- 3 *Super-resolution*: exploited residual training w/ low-res simulation data to ease the learning.

Fluid Simulation: a Modern Application of GNNs

Finite Volume Features, Global Geometry Representations, and Residual Training for Deep Learning-based CFD Simulation, ICML 2024 (spotlight)

- **Problem:** Predict velocity, pressure at every point in the domain using AI.
- **Motivation:**

- Numerical simulations of fluids are computationally expensive often requiring days.
- Recently, data-driven GNNs have been deployed as an efficient alternative. (e.g. MeshGraphNet from Deepmind)
- Can we improve the fidelity of data-driven GNNs?

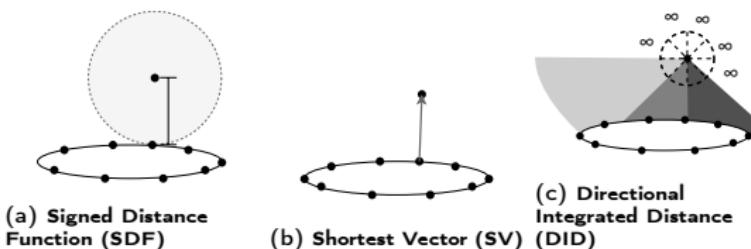


(a) CFD mesh with an airfoil body surrounded by different sizes of cells. (b) Illustration of cell characteristics, namely cell centroids, x_i and x_j , face centroid, $c_{f,ij}$, face area normal vector, $S_{f,ij}$, and cell volumes, V_i and V_j .

Contributions

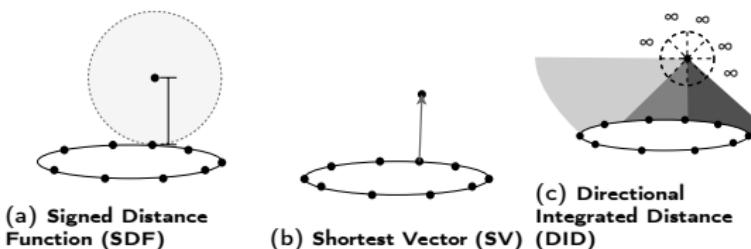
- 1 Geometric features to inform GNN nodes of long-range interactions.
- 2 Finite-volume features in the graph convolution layer.
- 3 *Super-resolution*: exploited residual training w/ low-res simulation data to ease the learning.

Impact of this work



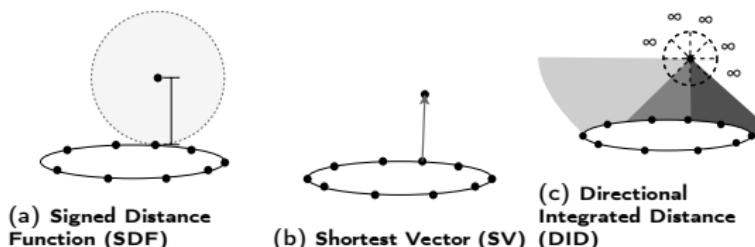
- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We showed that using cell characteristics, such as cell volume as node features, while face surface area, and face centroid as edge features improves the fidelity.
- 3 We were able to improve the fidelity of MeshGraphNet by 41% (and others).
- 4 Patent at UK IP office.

Impact of this work



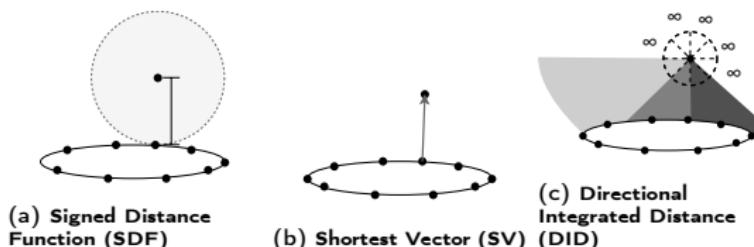
- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We showed that using cell characteristics, such as cell volume as node features, while face surface area, and face centroid as edge features improves the fidelity.
- 3 We were able to improve the fidelity of MeshGraphNet by 41% (and others).
- 4 Patent at UK IP office.

Impact of this work



- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We showed that using cell characteristics, such as cell volume as node features, while face surface area, and face centroid as edge features improves the fidelity.
- 3 We were able to improve the fidelity of MeshGraphNet by 41% (and others).
- 4 Patent at UK IP office.

Impact of this work



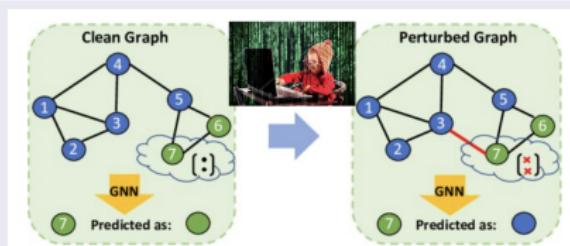
- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We showed that using cell characteristics, such as cell volume as node features, while face surface area, and face centroid as edge features improves the fidelity.
- 3 We were able to improve the fidelity of MeshGraphNet by 41% (and others).
- 4 Patent at UK IP office.

Adversarial attack: a Modern Challenge to GNNs

When Witnesses Defend: A Witness Graph Topological Layer for Adversarial Graph Learning. (AAAI'25)

Adversarial attack on Graph learning algorithms.

Attacker misleads a learning algorithm (e.g. GNN) into making incorrect predictions or classifications by deliberately perturbing a small number of edges (e.g. remove/add edges) or node features.



Adversarial perturbation (around target node 7) causes misclassification.

Contributions

- 1 We introduced a novel topological adversarial defense, namely, the *Witness Graph Topological Layer (WGTL)*.
- 2 WGTL integrates not only local but also global higher-order graph characteristics and controls their potential defense role via a topological regularizer.

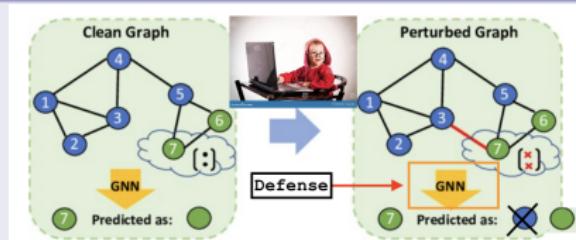
Adversarial attack: a Modern Challenge to GNNs

When Witnesses Defend: A Witness Graph Topological Layer for Adversarial Graph Learning. (AAAI'25)

Adversarial attack on Graph learning algorithms.

Attacker misleads a learning algorithm (e.g. GNN) into making incorrect predictions or classifications by deliberately perturbing a small number of edges (e.g. remove/add edges) or node features.

Problem



Design a defense algorithm that mitigates the effect of adversarial attack

Contributions

- 1 We introduced a novel topological adversarial defense, namely, the *Witness Graph Topological Layer (WGTL)*.
- 2 WGTL integrates not only local but also global higher-order graph characteristics and controls their potential defense role via a topological regularizer.



Topological Features

Topological Features

diagram B —————

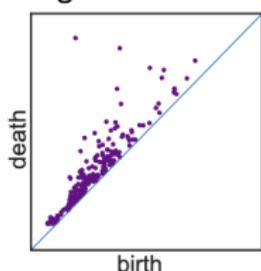
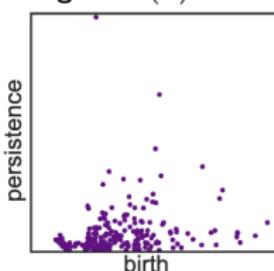


diagram $T(B)$ —————



A persistence diagram is transformed using function $T : (x, y) \rightarrow (0, y - x)$.

Why Topological features?

Stability theorem: Small ($\leq \epsilon$) change in the data (graph/points) only result in small ($\leq \epsilon$) changes in the persistence diagram.

Topological Features

diagram B —————

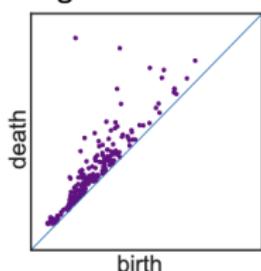
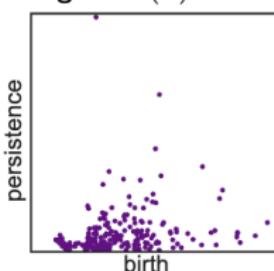


diagram $T(B)$ —————



A persistence diagram is transformed using function $T : (x, y) \rightarrow (0, y - x)$.

Why Topological features?

Stability theorem: Small ($\leq \epsilon$) change in the data (graph/points) only result in small ($\leq \epsilon$) changes in the persistence diagram.

Impact of this work

- ➊ This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- ➋ Effective against several types of attacks, for instance,
 - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
 - Global poisoning attack (Graybox, instead of targetting specific neighborhood modify whichever edges maximizes attackers objective)
 - Adaptive attacks (White-box, assumes that the model architecture, parameters and defense mechanisms are known to the attacker)
 - Node feature attack
- ➌ WGTL improves existing defenses such as Pro-GNN, GNNGuard, and SimP-GCN respectively by 5%, 15%, and 5%.

Impact of this work

- ➊ This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- ➋ Effective against several types of attacks, for instance,
 - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
 - Global poisoning attack (Graybox, instead of targetting specific neighborhood modify whichever edges maximizes attackers objective)
 - Adaptive attacks (White-box, assumes that the model architecture, parameters and defense mechanisms are known to the attacker)
 - Node feature attack
- ➌ WGTL improves existing defenses such as Pro-GNN, GNNGuard, and SimP-GCN respectively by 5%, 15%, and 5%.

Impact of this work

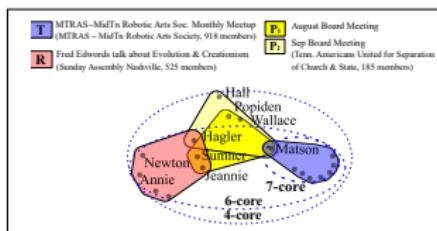
- ➊ This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- ➋ Effective against several types of attacks, for instance,
 - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
 - Global poisoning attack (Graybox, instead of targetting specific neighborhood modify whichever edges maximizes attackers objective)
 - Adaptive attacks (White-box, assumes that the model architecture, parameters and defense mechanisms are known to the attacker)
 - Node feature attack
- ➌ WGTL improves existing defenses such as Pro-GNN, GNNGuard, and SimP-GCN respectively by 5%, 15%, and 5%.

Hypergraphs: a modern paradigm in higher-order graph data mining

Neighborhood based hypergraph core decomposition (pVLDB 2023).

Neighborhood-based core decomposition

Decomposition of a hypergraph into nested, maximal subhypergraphs/cores such that all nodes in the k -core have at least k neighbors in that subhypergraph.



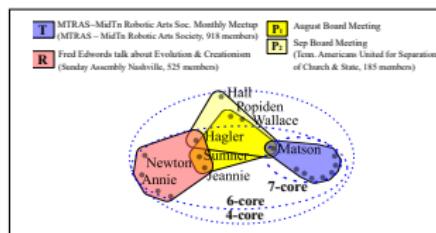
The set of events $H = \{T, R, P_1, P_2\}$ forms a hypergraph. 6-core $\Rightarrow \{T, R\}$, 7-core $\Rightarrow \{T\}$

Hypergraphs: a modern paradigm in higher-order graph data mining

Neighborhood based hypergraph core decomposition (pVLDB 2023).

Neighborhood-based core decomposition

Decomposition of a hypergraph into nested, maximal subhypergraphs/cores such that all nodes in the k -core have at least k neighbors in that subhypergraph.



The set of events $H = \{T, R, P_1, P_2\}$ forms a hypergraph. 6-core => $\{T, R\}$, 7-core => $\{T\}$

Contributions

- We introduced this novel core decomposition.
- We proposed an efficient local algorithm that scales to million-node hypergraph.
- **Applications:**
 - **Densest subhypergraph extraction.** Our novel volume-densest subhypergraphs capture important meetup events.
 - **Diffusion intervention.** Our decomposition is more effective than other graph-based decompositions in intervening diffusion e.g. epidemic spread.

Impact of this work

- ① Fastest hypergraph core-decomposition algorithm to date: Can decompose ArnetMiner hypergraph with 27M nodes and 17M hyperedges in-memory within 91 seconds.
- ② Laid the foundation for more recent core decomposition methods such as (k,g) -core [CIKM'23], (k,t) -hypercore [ECML-PKDD'23], and Dual-Layer Hierarchy [CIKM'24].

(Brief) Future Plans

Problems I want to study and relevant impact areas:

① Cyber-security.

- **Privacy-preserving graph neural networks** that are robust against adversaries.
- **Adversarial robustness of Higher-order graph (hypergraph) learning.**
- Efficient **uncertainty estimation of GNNs**. We want increased trustworthiness of AI systems by providing transparency about model confidence. In cybersecurity, uncertainty estimation can help prioritize high-risk alerts for further investigation.

② Human-centered AI.

- **Inconsistencies and bias in reasoning ability of LLMs.** Can we develop benchmark datasets to assess the deductive reasoning ability of LLMs?
- How to model rich datasets e.g. those arising from **multiple modalities and multiple sources** using a **Knowledge Hypergraph**. How to handle complex queries on such knowledge hypergraph?

③ Physical science.

- **Explainability and reliability of ML models** (particularly GNNs) in the physics domain e.g. fluid dynamics, additive manufacturing etc.

(Brief) Future Plans

Problems I want to study and relevant impact areas:

① Cyber-security.

- **Privacy-preserving graph neural networks** that are robust against adversaries.
- **Adversarial robustness of Higher-order graph (hypergraph) learning.**
- Efficient **uncertainty estimation of GNNs**. We want increased trustworthiness of AI systems by providing transparency about model confidence. In cybersecurity, uncertainty estimation can help prioritize high-risk alerts for further investigation.

② Human-centered AI.

- **Inconsistencies and bias in reasoning ability of LLMs.** Can we develop benchmark datasets to assess the deductive reasoning ability of LLMs?
- How to model rich datasets e.g. those arising from **multiple modalities and multiple sources** using a **Knowledge Hypergraph**. How to handle complex queries on such knowledge hypergraph?

③ Physical science.

- Explainability and reliability of ML models (particularly GNNs) in the physics domain e.g. fluid dynamics, additive manufacturing etc.

(Brief) Future Plans

Problems I want to study and relevant impact areas:

① Cyber-security.

- **Privacy-preserving graph neural networks** that are robust against adversaries.
- **Adversarial robustness of Higher-order graph (hypergraph) learning.**
- Efficient **uncertainty estimation of GNNs**. We want increased trustworthiness of AI systems by providing transparency about model confidence. In cybersecurity, uncertainty estimation can help prioritize high-risk alerts for further investigation.

② Human-centered AI.

- **Inconsistencies and bias in reasoning ability of LLMs.** Can we develop benchmark datasets to assess the deductive reasoning ability of LLMs?
- How to model rich datasets e.g. those arising from **multiple modalities and multiple sources** using a **Knowledge Hypergraph**. How to handle complex queries on such knowledge hypergraph?

③ Physical science.

- **Explainability and reliability** of ML models (particularly GNNs) in the **physics domain** e.g. fluid dynamics, additive manufacturing etc.

Year 1-2 (Foundations and Early Collaborations)

Develop private+robust GNNs and series of collaborations with Virginia Tech, and GVSU colleagues.

NSF SatC 2.0
NSF CISE/CNS
NIST ITL MSE
NIH

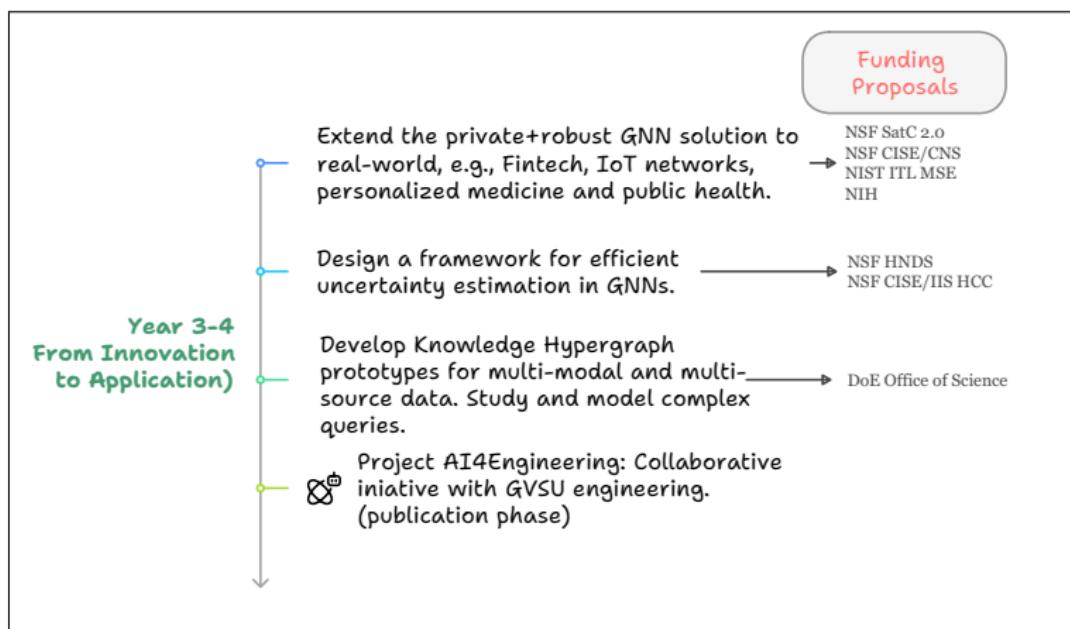
Develop and release benchmark dataset for assessing deductive reasoning in LLMs.

NSF HNDS
NSF CISE/IIS HCC

Project AI4Engineering: Collaborative initiative with GVSU engineering. (Planning and Funding phase)

DoE Office of Science

Mid term (years 3-4)



Long term (year 5)

Funding
Proposals

- Year 5
(Expansion and
Long-term impact)**
- Implement Uncertainty Estimation methods in collaboration with cybersecurity industry partners.
 - Secure a multi-institutional grant for Graph Learning for cyber-security and Uncertainty estimation in GNNs. Organize and host a workshop. → To be decided
 - Deploy Knowledge Hypergraph models in interdisciplinary applications (e.g., healthcare).
 - Project AI4Engineering: Collaborative initiative with GVSU engineering. (Extend the scope to other disciplines) → To be decided

Interdisciplinary Works and Collaborations

Interdisciplinary collaboration

- **Rolls-Royce plc.**, Dept. of Mechanical & Aerospace Engineering (NTU)
 - 1 Patent @UK IP office, 1 ICML paper.
 - Rolls-Royce has plans to internally adopt our solution.

Collaborations across the world

- *North-America:*
 - Virginia Tech
 - Purdue University
 - University of California Riverside
 - Pacific Northwest National Lab (PNNL)
- *Europe:*
 - Aalborg University, Denmark
 - Inria @ Univ. of Lile, France
 - University of Vienna, Austria
 - CENTAI, Italy
 - Max Planck Institute, Germany
- *Asia:* NUS, NTU (Singapore), BUET (Bangladesh).

Interdisciplinary Works and Collaborations

Interdisciplinary collaboration

- **Rolls-Royce plc.**, Dept. of Mechanical & Aerospace Engineering (NTU)
 - 1 Patent @UK IP office, 1 ICML paper.
 - Rolls-Royce has plans to internally adopt our solution.

Collaborations across the world

- *North-America:*
 - **Virginia Tech**
 - **Purdue University**
 - **University of California Riverside**
 - **Pacific Northwest National Lab (PNNL)**
- *Europe:*
 - Aalborg University, **Denmark**
 - Inria @ Univ. of Lile, **France**
 - University of Vienna, **Austria**
 - CENTAI, **Italy**
 - Max Planck Institute, **Germany**
- *Asia:* NUS, NTU (Singapore), BUET (Bangladesh).

Mentoring

I have had the pleasure to mentor and collaborate with the following young researchers:

- **PhD students:**

- Siddhartha Shankar Das, PhD student@Purdue University (2024-) (To join PNNL as post-doc next summer)
- Bishwamittra Ghosh, PhD student@NUS (2021-2023)(Now post-doc at Max Planck Institute, Germany)
- Ehsan B. Mobaraki, PhD student@AAU (2023-2024)
- Sarah Hasan, PhD student@AAU (2023-2024)
- Loh Sher En Jessica, PhD student@NTU (2021-2024)
- Debabrata Mahapatra, PhD student@NUS (2019)

- **Undergraduate student:**

- Arpit Kumar Rai, IIT Kanpur, Intern @NTU (2021-2022) (Now Software Engineer @ Glean, Palo Alto, California)

Thank You

Q&A

Slides

