

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

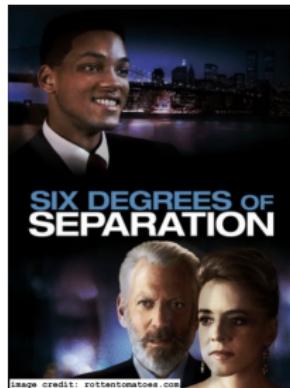
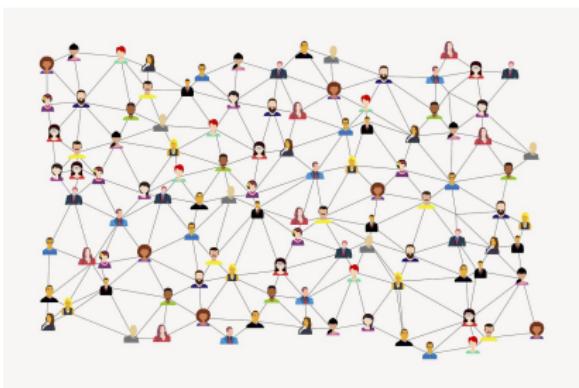
4 Future Plans

5 Interdisciplinary Works and Collaborations

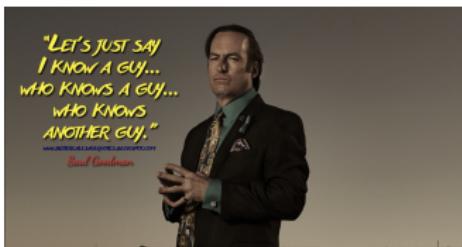
6 Mentoring

7 Q&A

Graphs before..



Graphs before..



Graphs now..

Well, they are everywhere

Biology

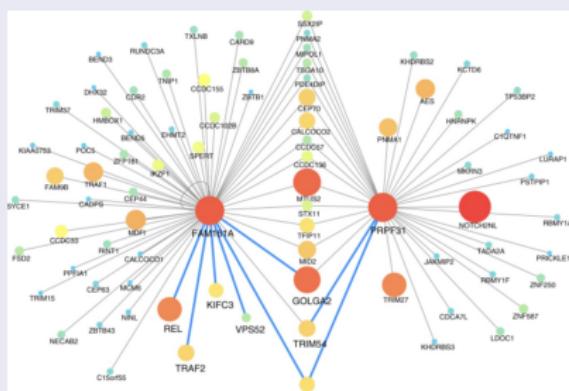
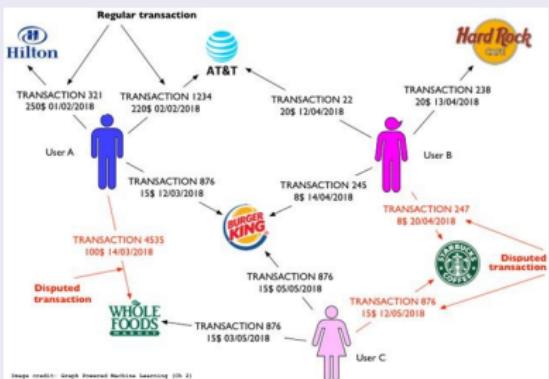


Image credit: Kovács, István A., et al. "Network-based prediction of protein interactions." *Nature comm.*, 2019

Graphs now..

Well, they are everywhere

Finance

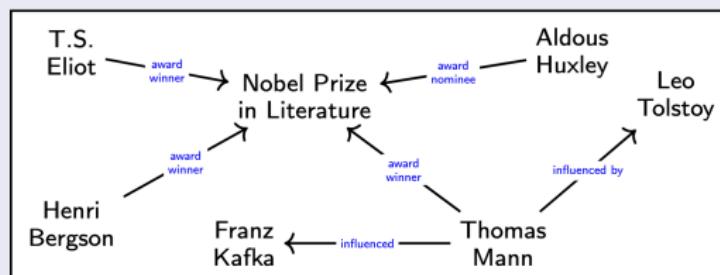


Financial transactions

Graphs now..

Well, they are everywhere

Human centered computing



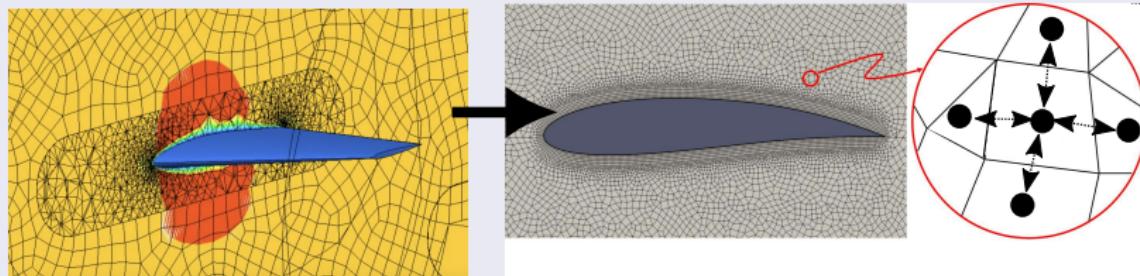
Knowledge graphs for explainable AI

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



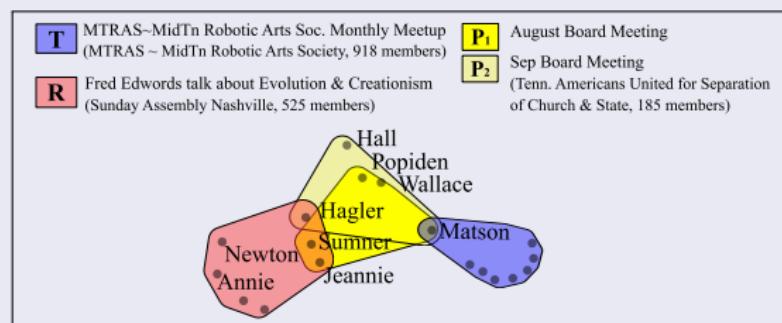
Fluid simulation at cross-section of an airplane wing (aka airfoil) and a close-up view

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

Graphs now..

Well, they are everywhere

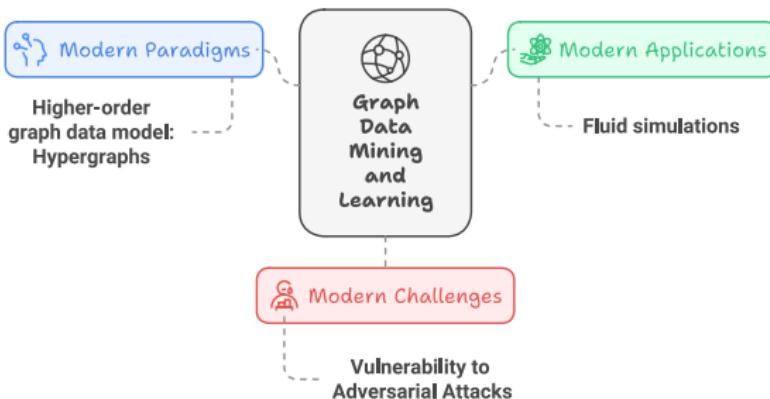
Public Health and Epidemiology



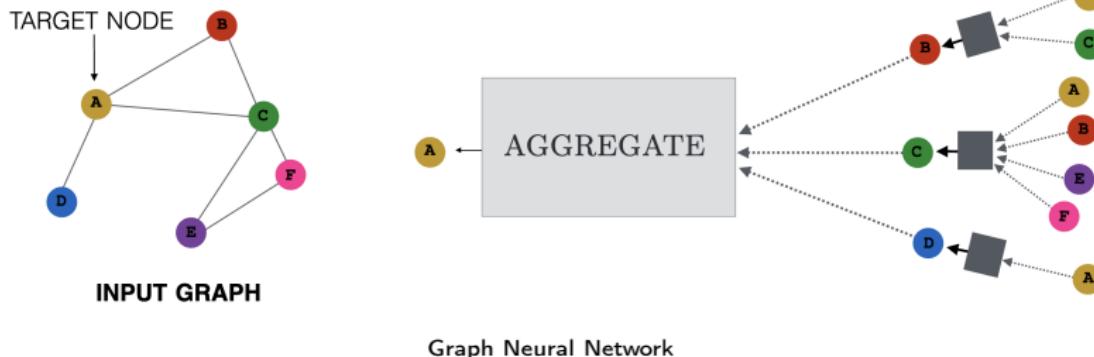
Modern paradigm: Modeling higher-order interactions as Hypergraphs.

My Research: Overview

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



My Research: Learning on Graph

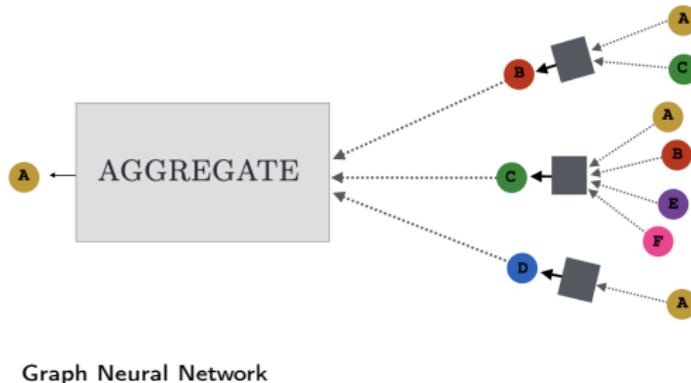
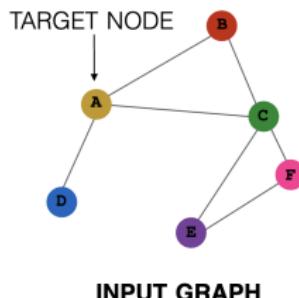


When GNNs (applications):

- 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

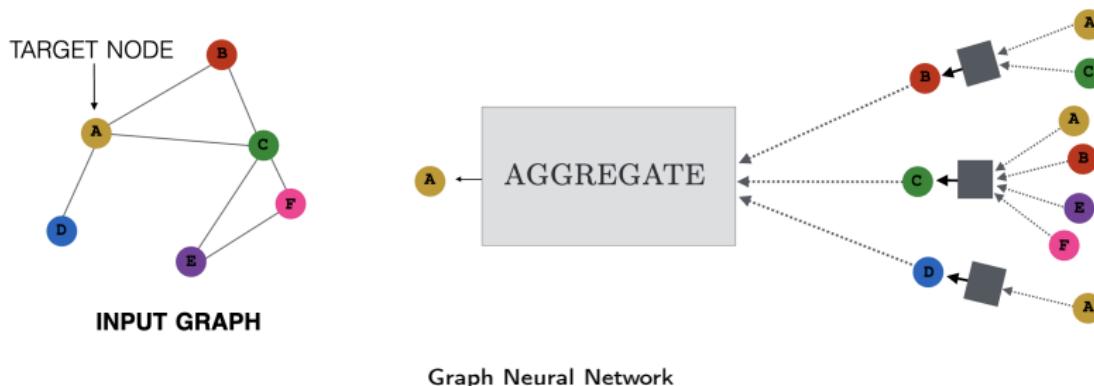


When GNNs (applications):

- 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

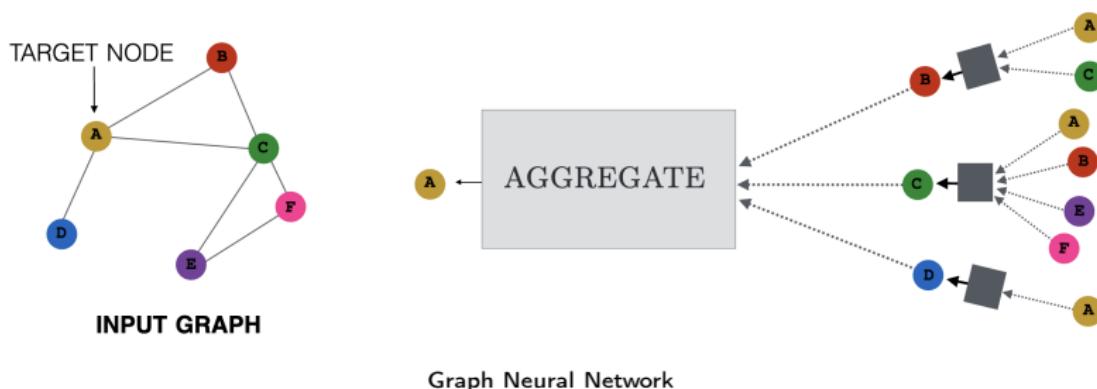


When GNNs (applications):

- 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.

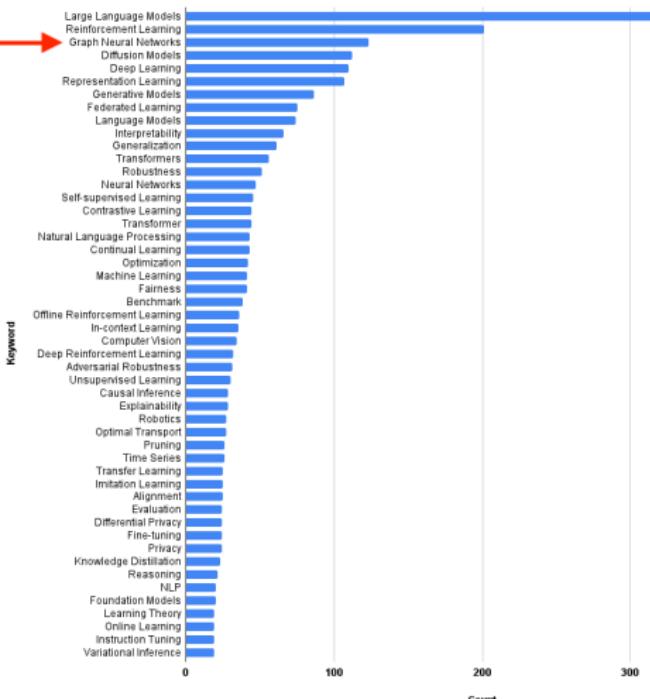
My Research: Why GNNs?

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



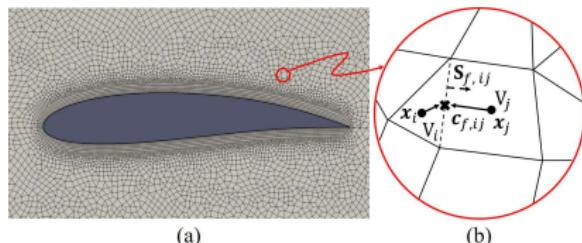
Fluid Simulation: a Modern Application of GNNs

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.

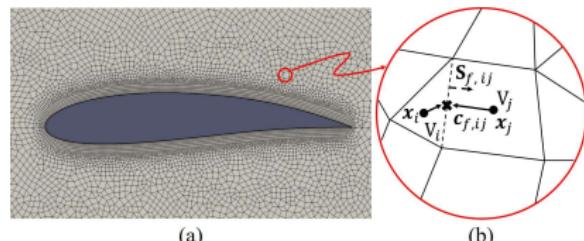
Fluid Simulation: a Modern Application of GNNs

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?

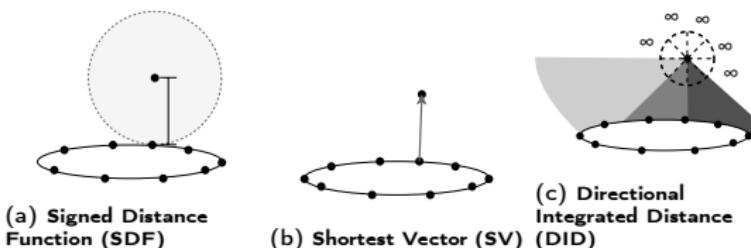


(a) CFD mesh with an airfoil body surrounded by different sizes of cells. (b) Illustration of cell characteristics, namely cell centroids, x_i and x_j , face centroid, $c_{f,ij}$, face area normal vector, $S_{f,ij}$, and cell volumes, V_i and V_j .

Contributions

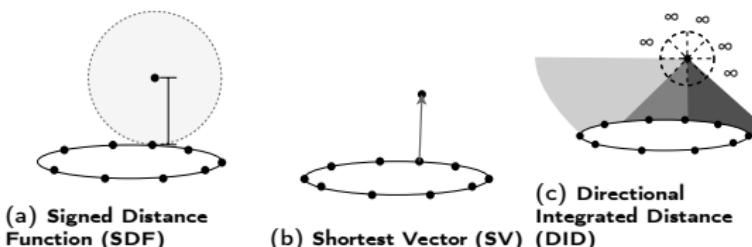
- ① Geometric features to inform GNN nodes of long-range interactions.
- ② Finite-volume features in the graph convolution layer.
- ③ *Super-resolution*: exploited residual training w/ low-res simulation data to ease the learning.

Impact of this work



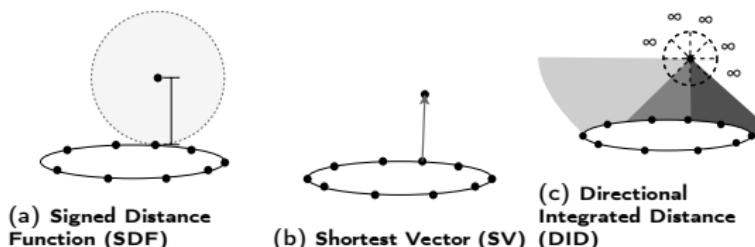
- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We 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



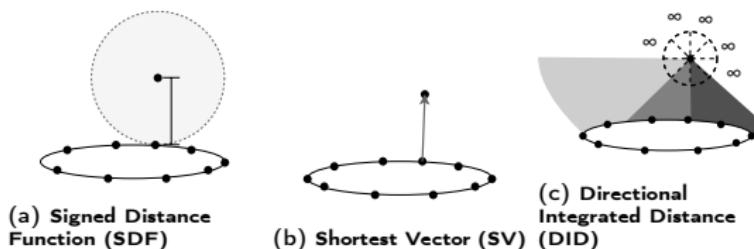
- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We 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



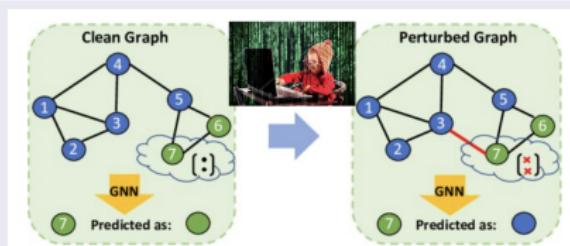
- 1 Existing geometric features (SDF) only indicate the presence of the closest boundary point somewhere along the circle's circumference. SV provides both distance and direction from the nearest boundary point. DID gives the average distance of all boundary points within a given angle range.
- 2 We 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.

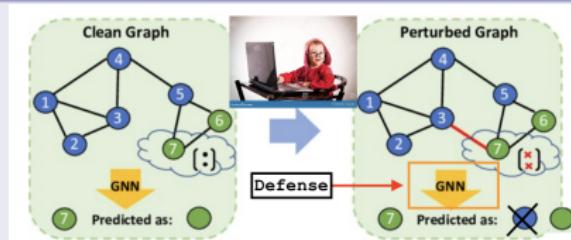
Adversarial attack: a Modern Challenge to GNNs

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 —————

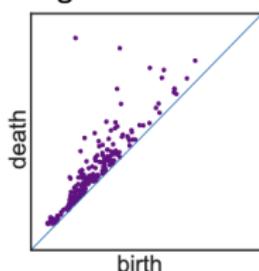
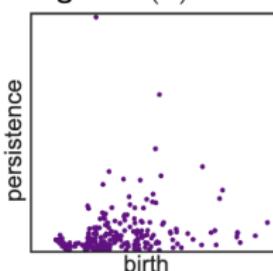


diagram $T(B)$ —————



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 —————

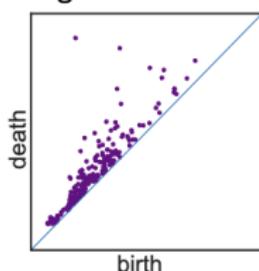
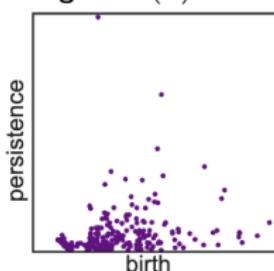


diagram $T(B)$ —————



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

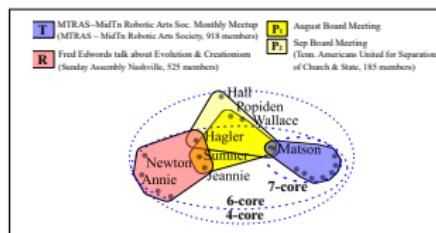
- ➊ This is the first work that shows that topological features can make GNNs robust against adversarial perturbations.
- ➋ Effective against several types of attacks, for instance,
 - Targetted poisoning attack (Graybox, modify the neighbors of a target node and their features)
 - Global poisoning attack (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 => $\{T, R\}$, 7-core => $\{T\}$

Contributions

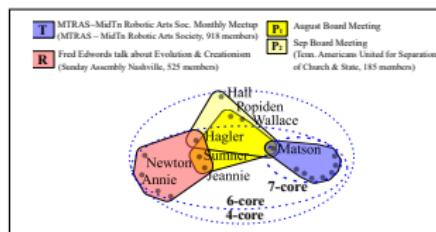
- We introduced this novel core decomposition.
- We proposed an efficient local algorithm that scales to million-node hypergraph.
- Applications:
 - Densest subhypergraph extraction. Our novel volume-densest subhypergraphs capture important meetup events.
 - Diffusion intervention. Our decomposition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

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.

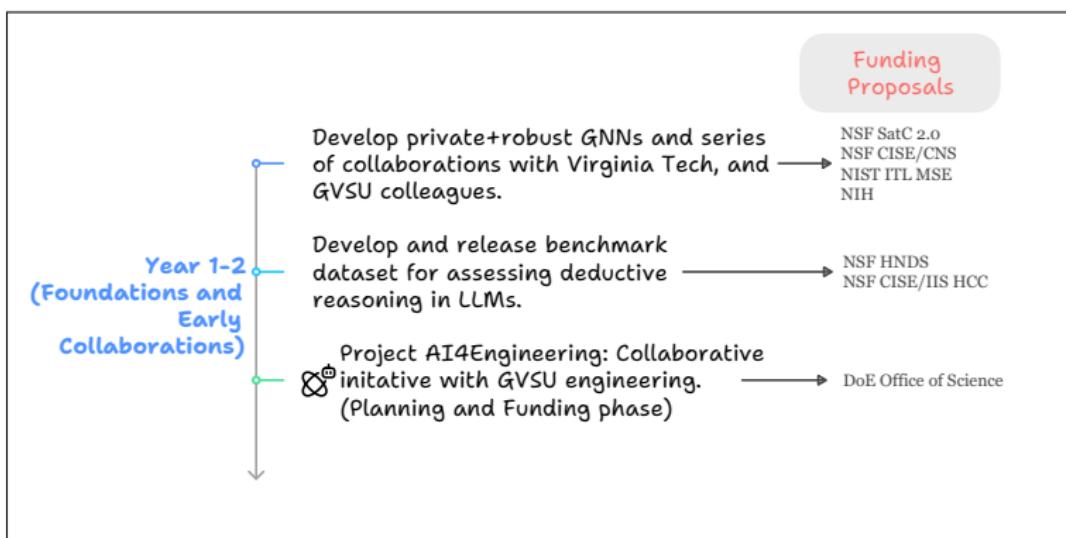
② Human-centered AI.

- **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?

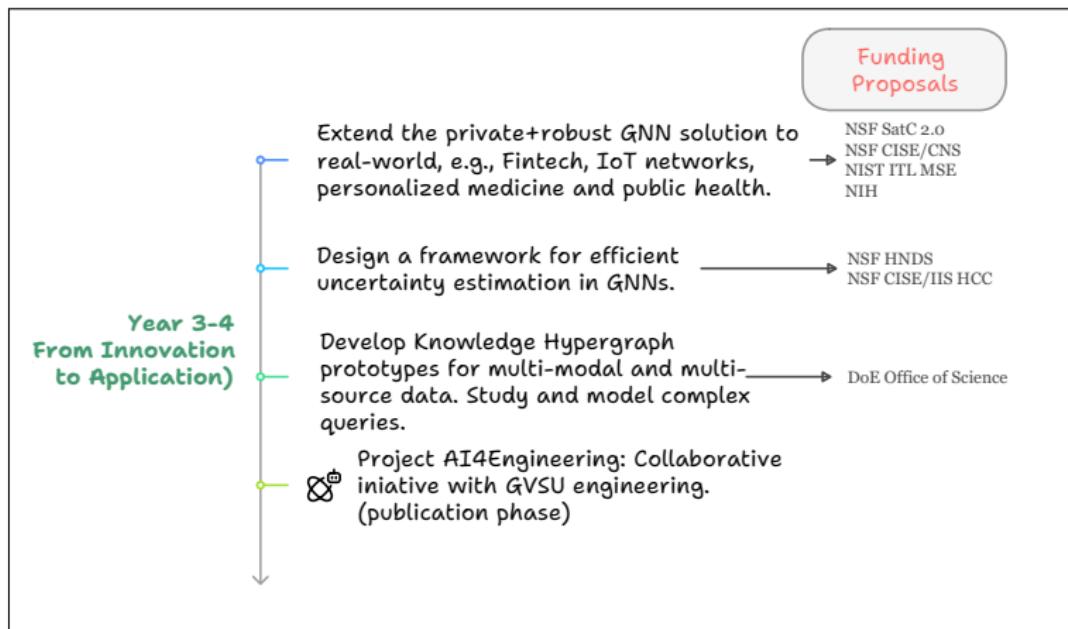
③ Physical science.

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

Short term (years 1-2)



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)
- Bishwamitra 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



Supplementary Slides





CFD+AI



Traditional (Numerical) Fluid Simulation

Step 1. Modeling the laws of fluid flow (Physics) with Equations.

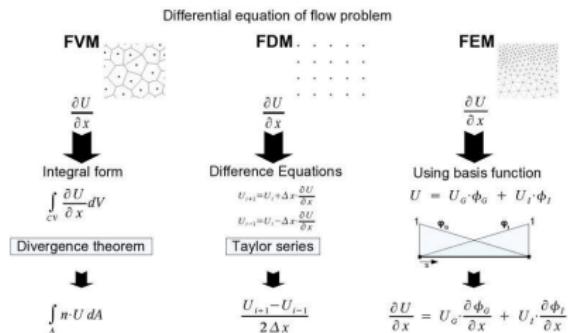
$$\begin{aligned}\frac{\partial u}{\partial t} + u \frac{\partial u}{\partial x} + v \frac{\partial u}{\partial y} &= -\frac{1}{Q} \frac{\partial p}{\partial x} + v \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) \\ \frac{\partial v}{\partial t} + u \frac{\partial v}{\partial x} + v \frac{\partial v}{\partial y} &= -\frac{1}{Q} \frac{\partial p}{\partial y} + v \left(\frac{\partial^2 v}{\partial x^2} + \frac{\partial^2 v}{\partial y^2} \right) \\ \frac{\partial^2 p}{\partial x^2} + \frac{\partial^2 p}{\partial y^2} &= -Q \left(\frac{\partial u}{\partial x} \frac{\partial u}{\partial x} + 2 \frac{\partial u}{\partial y} \frac{\partial v}{\partial x} + \frac{\partial v}{\partial y} \frac{\partial v}{\partial y} \right)\end{aligned}$$

2D Navier-Stokes Equation

Step 2. Conditions: Prescribe the fluid properties (e.g. density, viscosity), boundary conditions of the domain, and initial conditions for the flow.

Traditional (Numerical) Fluid Simulation (Contd.)

Step 3. Discretization: To solve these equations, the fluid domain is divided into a large number of small regions (grids or control volumes), and the equations are solved numerically for each region.



Various methods of approximating the terms in governing PDEs.

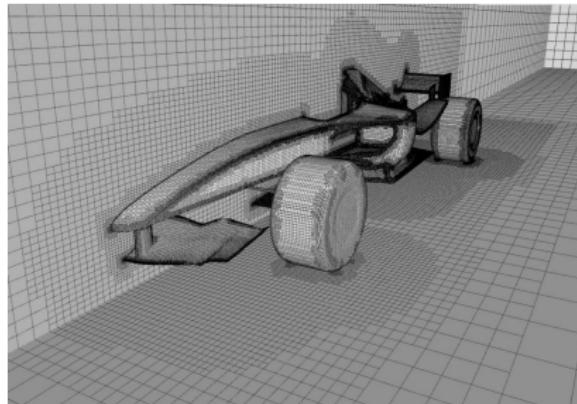
For instance, FDM linearizes the first equation at step 1

$$\frac{u_{i,j}^{n+1} - u_{i,j}^n}{\Delta t} + u_{i,j}^n \frac{u_{i,j}^n - u_{i-1,j}^n}{\Delta x} + v_{i,j}^n \frac{u_{i,j}^n - u_{i,j-1}^n}{\Delta y} = \\ - \frac{1}{\rho} \frac{P_{i+1,j} - P_{i-1,j}}{2\Delta x} + v \left(\frac{u_{i+1,j}^n - 2u_{i,j}^n + u_{i-1,j}^n}{\Delta x^2} + \frac{u_{i,j+1}^n - 2u_{i,j}^n + u_{i,j-1}^n}{\Delta y^2} \right)$$

Step 4. Solve the system of equations using numerical solvers (e.g. Gauss-Siedel).

Why Numerical Solvers are expensive?

- ① Ensuring **numerical stability and accuracy** often requires fine grids and small time steps, which leads to solving millions or even billions of coupled equations iteratively.
- ② Real-world objects are complex with many **complicated parts design**.



Finite Volume Features

Given node feature $h_i \in \mathcal{R}^{d_{in}}$ at the i^{th} node and its spatial location $x_i \in \mathbb{R}^t$, a convolution using FVF is defined as

$$\bar{h}_i(U, b, N_i) = \sum_{j \in N_i} \text{ReLU}(U^T(q_{ij}) + b) \odot (h_j \oplus p_j), \quad (1)$$

where $U \in \mathbb{R}^{3t \times (d_{in}+1)}$ and $b \in \mathbb{R}^{d_{in}+1}$ are trainable parameters, t is the spatial dimension of the CFD simulation, and N_i is an index set indicating the neighbourhood of the node i and \odot is the element-wise multiplication.

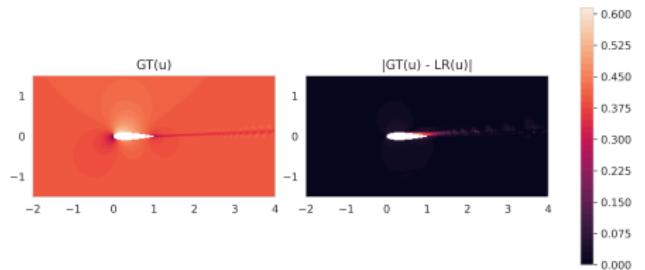
Finite Volume Features(Contd.)

In our model, the node attributes are its associated cell volume, denoted as $p_j = V_j \in \mathbb{R}^1$, and the edge attributes are its associated face area normal vector and the relative spatial location of its face centroid to the nodes, denoted as

$$q_{ij} = \mathbf{S}_{f,ij} \oplus (\mathbf{c}_{f,ij} - \mathbf{x}_i) \oplus (\mathbf{c}_{f,ij} - \mathbf{x}_j) \in \mathbb{R}^{3t}$$

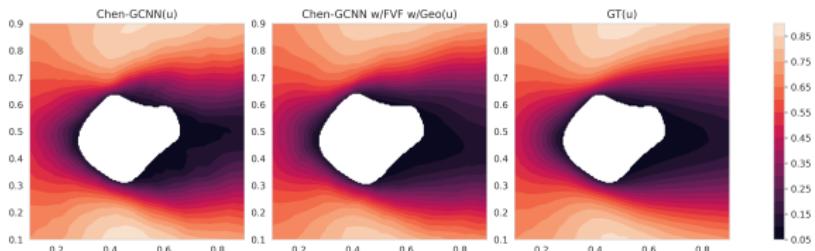
Residual training

- Train the network to predict the residual field ($GT - LR$), where LR is the upsampled low-resolution field, instead of the original field GT itself.
- Much of the residual field will be close to zero, which
 - eases the learning,
 - helps the model to focus on areas where the LR fields tend to be inaccurate



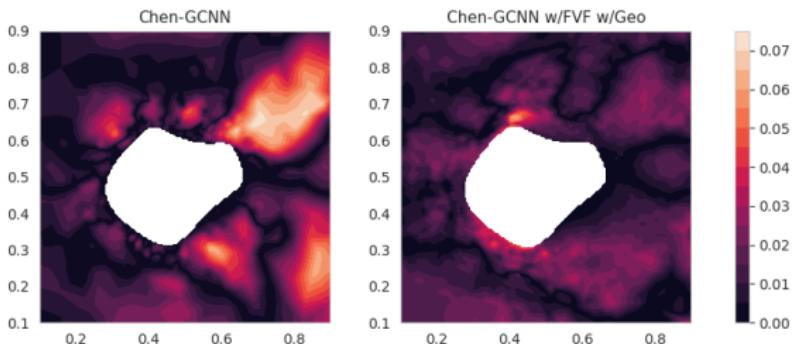
Prediction on an airfoil shows that the velocity x-component, u remains invariant in majority regions (left), hence the absolute value of residual ($|GT(u) - \text{Upsample}(LR(u))|$) is close to 0 in those regions (right).

Some visualisations



(a) Velocity x-component (u) visualisation.

$$|\hat{U} - U|$$



(b) Absolute error of predictions.

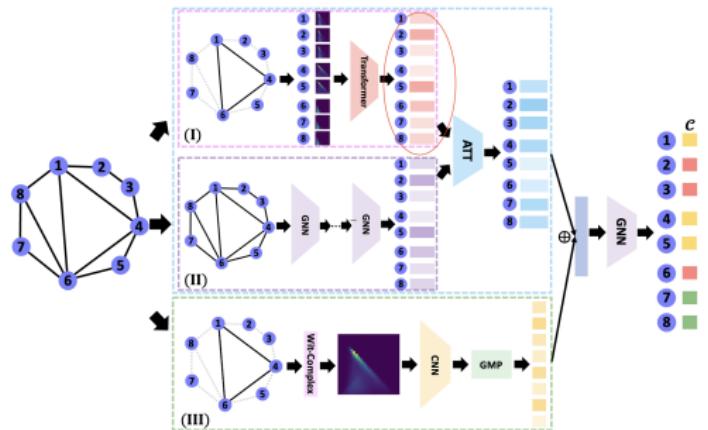
Chen-GCNN w/ FVF w/ Geo has a comparatively smaller high-error ([0.050, 0.075]) region than the baseline. Darkest shade = low-error region.

Efficiency of Geometric feature computations.

Average Computation Time (s)		
	Coarse-AirfRANS	AirfRANS
SV	0.001	0.018
DID	0.031	4.276
FVF	0.001	0.032

Computational time of SV, DID and FVF at two mesh resolutions. Airfoils in Coarse-AirfRANS has ~ 20 times fewer nodes than AirfRANS.

Adversarial robustness



Architecture of Witness Graph Topological Layer.

1 Local Topology Encoding: Encodes local topological features of every node. (Z_{T_L})

2 Node Representation Learning. Learns node representations using any backbone GNN. (Z_G)

3 Global Topology Encoding. Encodes topological feature of the entire graph. (Z_{T_G})

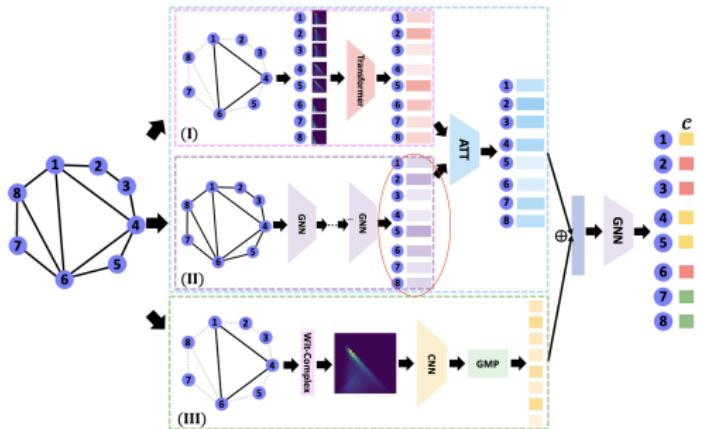
4 Aggregated Topological Encoding. Encodes local and global topological priors. (Z_{WGTL})

$$z = [Z_{T_L}, Z_G]$$

Attention coefficients, $\alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$

Additive attention, $Z_{AGG} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$

$$Z_{WGTL} = Z_{AGG} Z_{T_G}$$



Architecture of Witness Graph Topological Layer.

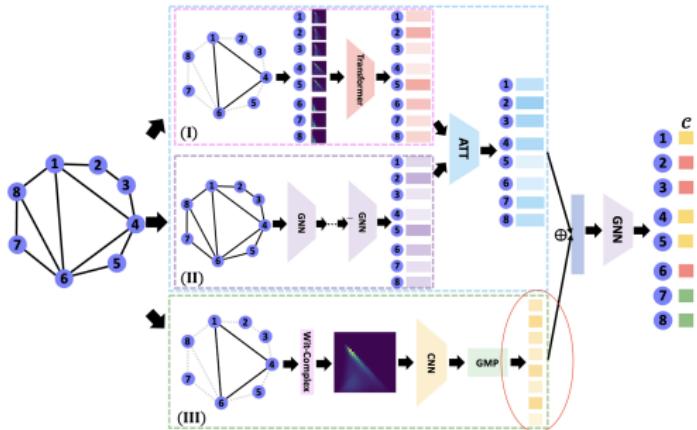
- ➊ Local Topology Encoding: Encodes local topological features of every node. (Z_{T_L})
- ➋ Node Representation Learning: Learns node representations using any backbone GNN. (Z_G)
- ➌ Global Topology Encoding: Encodes topological feature of the entire graph. (Z_{T_G})
- ➍ Aggregated Topological Encoding: Encodes local and global topological priors. (Z_{WGTL})

$$z = [Z_{T_L}, Z_G]$$

Attention coefficients, $\alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$

Additive attention, $Z_{\text{AGG}} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$

$$Z_{\text{WGTL}} = Z_{\text{AGG}} Z_{T_G}$$



Architecture of Witness Graph Topological Layer.

- 1 Local Topology Encoding: Encodes local topological features of every node. (Z_{T_L})
- 2 Node Representation Learning: Learns node representations using any backbone GNN. (Z_G)
- 3 Global Topology Encoding: Encodes topological feature of the entire graph. (Z_{T_G})

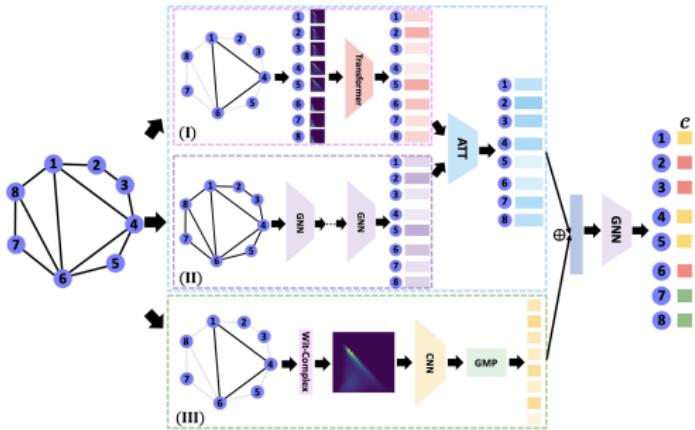
- 4 Aggregated Topological Encoding: Encodes local and global topological priors. (Z_{WGTL})

$$z = [Z_{T_L}, Z_G]$$

Attention coefficients, $\alpha_i = \text{Softmax}(W_2 \cdot \tanh(W_1 z_i + b_1))$

Additive attention, $Z_{\text{ADD}} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$

$$Z_{WGTL} = Z_{\text{ADD}} Z_{T_G}$$



Architecture of Witness Graph Topological Layer.

- I Local Topology Encoding: Encodes local topological features of every node. (Z_{T_L})
- II Node Representation Learning: Learns node representations using any backbone GNN. (Z_G)
- III Global Topology Encoding: Encodes topological feature of the entire graph. (Z_{T_G})
- IV Aggregated Topological Encoding: Encodes local and global topological priors. (Z_{WGTL})

$$z = [Z_{T_L}, Z_G]$$

Attention coefficients, $\alpha_i = \text{Softmax}(\mathbf{W}_2 \cdot \tanh(\mathbf{W}_1 z_i + \mathbf{b}_1))$

Additive attention, $Z_{AGG} = \alpha_1 \times Z_{T_L} + \alpha_2 \times Z_G$

$$Z_{WGTL} = Z_{AGG} Z_{T_G}$$

WGTL: Topological Regularizer

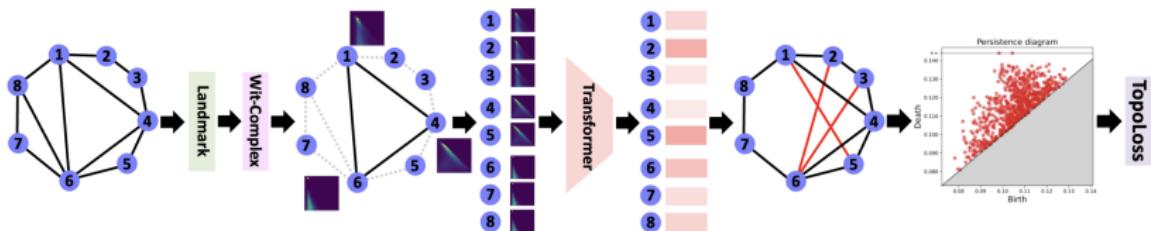


Illustration of Witness Complex-based topological regularizer L_{topo} .

$$L_{topo}(T(\mathcal{G})) \triangleq \sum_{i=1}^m (d_i - b_i)^2 \left(\frac{d_i + b_i}{2} \right)^2, \quad (2)$$

- A localized attack (perturbing certain nodes or edges) appears as topological noise in the final persistent diagram, and exhibit lower persistence.
- And minimising L_{topo} forces the Transformer to learn local topology encodings (Z_{T_L}) which produces PD with small persistence, i.e., $(d_i - b_i)$.

WGTL: Topological Regularizer

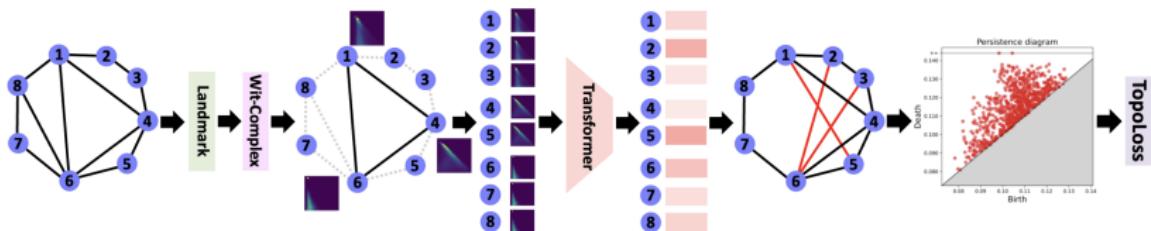


Illustration of Witness Complex-based topological regularizer L_{topo} .

$$L_{topo}(T(\mathcal{G})) \triangleq \sum_{i=1}^m (d_i - b_i)^2 \left(\frac{d_i + b_i}{2} \right)^2, \quad (2)$$

- A localized attack (perturbing certain nodes or edges) appears as topological noise in the final persistent diagram, and exhibit lower persistence.
- And minimising L_{topo} forces the Transformer to learn local topology encodings (Z_{T_L}) which produces PD with small persistence, i.e., $(d_i - b_i)$.

Results

Table 1: Comparison of performances (avg. accuracy \pm std.) with existing defenses under mettack.

Dataset	Models	Perturbation Rate		
		0%	5%	10%
Cora-ML	Pro-GNN	82.98 \pm 0.23	80.14 \pm 1.34	71.59 \pm 1.33
	Pro-GNN+WGTL	83.85\pm0.38	81.90\pm0.73	72.51\pm0.76
	GCN+GNNGuard	83.21 \pm 0.34	76.57 \pm 0.50	69.13 \pm 0.77
	GCN+GNNGuard+WGTL	*84.78\pm0.43	*83.23\pm0.82	*79.96\pm0.49
Citeseer	SimP-GCN	79.52 \pm 1.81	74.75 \pm 1.40	70.87 \pm 1.70
	SimP-GCN+WGTL	81.49\pm0.52	76.65\pm0.65	72.88\pm0.83
	ProGNN	72.34 \pm 0.99	68.96 \pm 0.67	67.36 \pm 1.12
	ProGNN+WGTL	72.83\pm0.94	71.85\pm0.74	70.70\pm0.57
Pubmed	GCN+GNNGuard	71.82 \pm 0.43	70.79 \pm 0.22	66.86 \pm 0.54
	GCN+GNNGuard+WGTL	73.37\pm0.63	72.57\pm0.17	66.93\pm0.21
	SimP-GCN	73.73 \pm 1.54	73.06 \pm 2.09	72.51 \pm 1.25
	SimP-GCN+WGTL	*74.32\pm0.19	*74.05\pm0.71	*73.09\pm0.50
Polblogs	Pro-GNN	87.33 \pm 0.18	87.25 \pm 0.09	87.20 \pm 0.12
	Pro-GNN + WGTL (ours)	87.90\pm0.30	*87.77\pm0.08	*87.67\pm0.22
	GCN+GNNGuard	83.63 \pm 0.08	79.02 \pm 0.14	76.58 \pm 0.16
	GCN+GNNGuard+WGTL	OOD	OOD	OOD
Polblogs	SimP-GCN	*88.11\pm0.10	86.98 \pm 0.19	86.30 \pm 0.28
	SimP-GCN+WGTL	OOD	OOD	OOD
	GCN+GNNGuard	95.03 \pm 0.25	73.25 \pm 0.16	72.76 \pm 0.75
	GCN+GNNGuard+WGTL	*96.22\pm0.25	*73.62\pm0.22	*73.72\pm1.00
Polblogs	SimP-GCN	89.78 \pm 6.47	65.75 \pm 5.03	61.53 \pm 6.41
	SimP-GCN+WGTL	94.56\pm0.24	69.78\pm4.10	69.55\pm4.42

Efficiency

Table 2: Efficiency of WGTL. All the times are in seconds.

Datasets/ (# Landmarks)	Landmark selection time	Local feat. comput. time	Global feat. comput. time
Cora-ML/124	0.01±0.01	0.12±0.03	5.11±0.13
Citeseer/105	0.01±0.01	0.16±0.02	5.23±1.22
Polblogs/61	0.01±0.00	0.07±0.01	4.64±0.2
Snap-patents/91	0.03±0.02	0.64±0.00	7.54±1.15
Pubmed/394	0.07±0.01	0.51±0.03	27.83±0.47
OGBN-arXiv/84	1.02 ± 0.00	12.79±0.31	83.04±2.19

Hypergraphs

Hypergraph

A **hypergraph** (V, E) consists of a set of nodes V and a collection of subsets of nodes E called *hyperedges*. Unlike edges in a graph, hyperedge may contain more than 2 nodes.

Examples: co-authorship in papers, event-participant relations in meet-ups, etc.

Neighbors

Pair of nodes that co-occur in a hyperedge are neighbours.

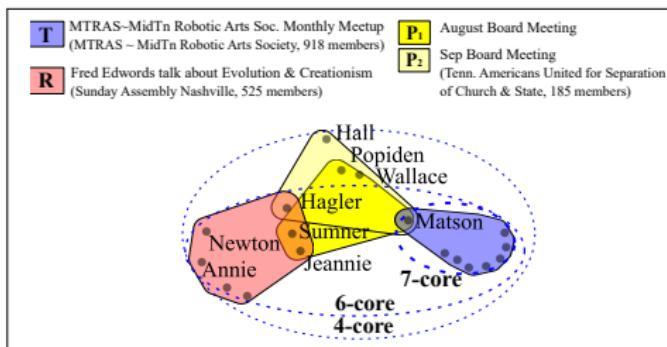


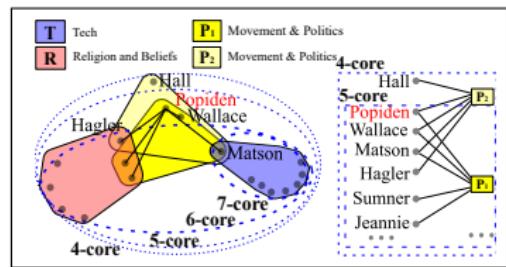
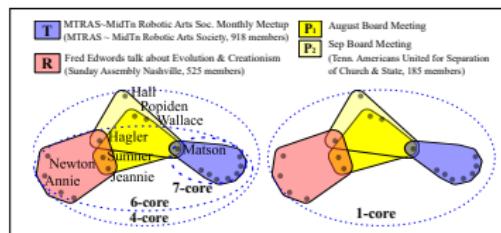
Figure 3: The set of events $H = \{T, R, P_1, P_2\}$ forms a hypergraph. *Annie* and *Newton* are neighbors. *Newton* has 6 neighbours.

Motivation

Limitations of existing methods.

Hypergraph Degree-based decomposition may not be informative

Reduced Hypergraph Reducing to Clique graph and bipartite graph and then applying graph core-decompositions produces non-intuitive results.



Challenges

Peeling paradigm In the classic peeling algorithm for graph, a node removal reduces its neighbors' degree by 1 (Linear time algorithm). However, in a neighborhood-based hypergraph core decomposition, its neighboring nodes # neighbors may reduce by more than 1 (Polynomial time).

Local algorithm paradigm Graph h -index reports incorrect neighborhood-based core.

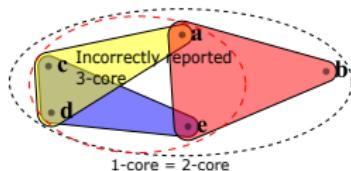


Figure 6: For any $n > 1$, the h -index of node a never reduces from $h_a^{(1)} = \mathcal{H}(2, 3, 3, 4) = 3$ to its correct core-number 2. Because a will always have at least 3 neighbors (c , d , and e) whose h -indices are at least 3. An incorrect 3-core reported.

Problem Statement

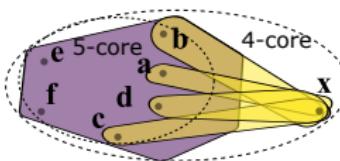
How to correctly and efficiently compute neighborhood-based hypergraph cores.

Naive Peeling algorithm: Peel

- At each iteration $k \in \{1, 2, \dots, |V|\}$,
 - Remove the node with $\# \text{ neighbors} \leq k$.
 - Report k as the core-number of the removed node.
 - Recompute the $\# \text{ neighbors}$ of neighboring nodes.
- Complexity: $\mathcal{O}(|V| \cdot d_{nbr} \cdot (d_{nbr} + d_{hpe}))$, here d_{nbr} (d_{hpe}) is the $\# \text{ neighbor}$ (degree) of the node with largest $\# \text{ neighbors}$ (degree).

Can we do better?.

Delay $\# \text{ neighbors}$ recomputation of nodes with core-number $> k$ based on lower-bound.



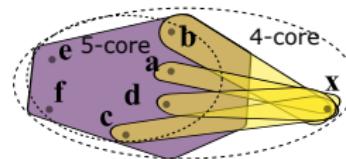
Node b 's $\# \text{ neighbors}$ is recomputed twice: (1) when x is peeled and later (2) when e is peeled. Can we delay the recomputation until e is peeled?

Better Peeling algorithm: E-Peel

- 1 Compute the core-number lower bound for all nodes.
- 2 At each iteration $k \in \{1, 2, \dots, |V|\}$,
 - 1 Remove the node with $\#$ neighbors $\leq k$.
 - 2 Report k as the core-number of the removed node.
 - 3 Recompute the $\#$ neighbors of a neighboring node v only if $k \geq LB(v)$.

$$LB(v) = \max \left(|e_m(v)| - 1, \min_{u \in V} |N(u)| \right)$$

Here $e_m(v)$ is the maximal cardinality hyperedge containing v



Node b 's $\#$ neighbors computation is delayed until e is peeled.

Best algorithm: Local core and optimisations

Input: Hypergraph $H = (V, E)$

Output: Core-number $c(v)$ for each node $v \in V$

for all $v \in V$ do

$\hat{h}_v^{(0)} = h_v^{(0)} \leftarrow |N(v)|$.

for all $n = 1, 2, \dots, \infty$ do

for all $v \in V$ do

$h_v^{(n)} \leftarrow \min \left(\mathcal{H}(\{\hat{h}_u^{(n-1)} : u \in N(v)\}), \hat{h}_v^{(n-1)} \right)$

for all $v \in V$ do

$c(v) \leftarrow \hat{h}_v^{(n)} \leftarrow \text{Core-correction } (v, h_v^{(n)}, H)$

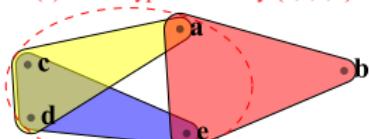
if $\forall v, \hat{h}_v^{(n)} == h_v^{(n)}$ then

Terminate Loop

Return c

Core correction:

$H^+(a) = \text{sub-hyp. induced by } \{a, e, c, d\}$



$h_a = h_e = h_c = h_d = 3$.

$h_b = 2$

$H^+(a) = H[\{u : h_u \geq h_a\}]$

Reduce h -index h_a by 1 until the #neighbors of a in $H^+(a) \geq h_a$: Node a 's corrected h -index = 2.

Hypergraph h -index of order n

The Hypergraph h -index of order n for node v , denoted as $\hat{h}_v^{(n)}$, is defined for any natural number $n \in \mathbb{N}$ by the following recurrence relation:

$$\hat{h}_v^{(n)} = \begin{cases} |N(v)| & n = 0 \\ h_v^{(n)} & n > 0 \wedge LCCSAT(h_v^{(n)}) \\ \max\{k \mid k < h_v^{(n)} \wedge LCCSAT(k)\} & n > 0 \wedge \neg LCCSAT(h_v^{(n)}) \text{ (Core-correction)} \end{cases} \quad (3)$$

Local core: Theoretical Guarantees

Hypergraph h-index has a limit

For any node $v \in V$ of a hypergraph $H = (V, E)$, the two sequences $(h_v^{(n)})$ and $(\hat{h}_v^{(n)})$ have the same limit: $\lim_{n \rightarrow \infty} h_v^{(n)} = \lim_{n \rightarrow \infty} \hat{h}_v^{(n)}$.

The limiting value is the core-number

If the local coreness-constraint is satisfied for all nodes $v \in V$ at the terminal iteration, the corrected h -index at the terminal iteration $\hat{h}_v^{(\infty)}$ satisfies: $\hat{h}_v^{(\infty)} = c(v)$.

Convergence time guarantee

Given a node $v \in N_i$ in a hypergraph H , it holds that $\forall n \geq i, \hat{h}_v^{(n)} = c(v)$. Here N_i is the i -th *neighborhood hierarchy*, which contains the set of nodes that have the minimum number of neighbors in $H[V']$, where $V' = V \setminus \cup_{0 \leq j < i} N_j$

Optimisations and parallelisation of Local core

- **Optimisations:** We have proposed 4 optimisations to make **Local-core** more efficient.
- **Parallisation:** We have proposed **Local-core(p)**, a shared-memory, data parallel programming adaptation of Local-core.
- **Generalised core model:** We have proposed a generalised hypergraph core model (*neighborhood, degree*)-core that simultaneously considers degree constraint and neighborhood constraint.

Efficiency evaluation

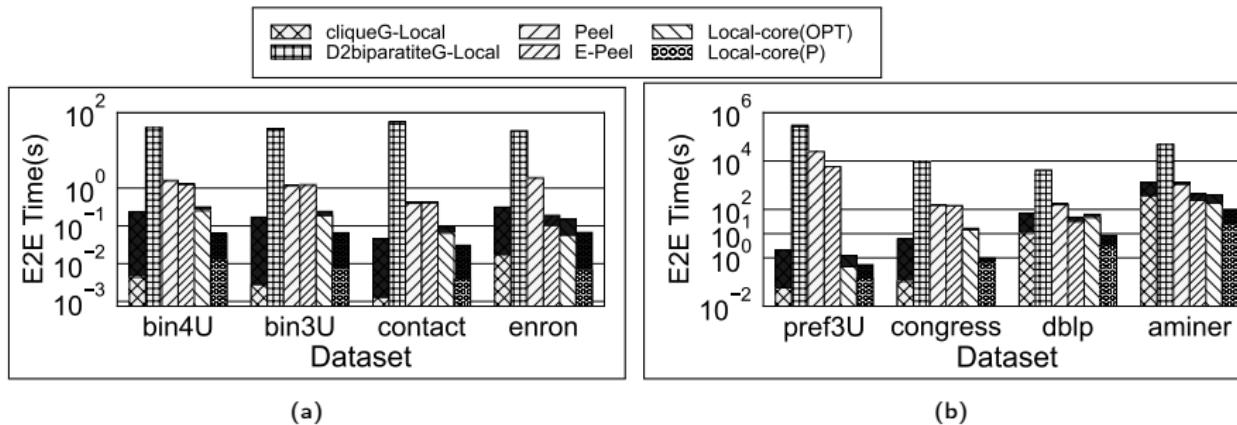


Figure 7: (a)-(b) End-to-end (E2E) running time of our algorithms: Peel, E-Peel, Local-core(OPT), Local-core(P) with 64 Threads vs. those of baselines: Clique-Graph-Local and Distance-2 Bipartite-Graph-Local. End-to-end (E2E) running time = data structure initialization time (shaded with dark-black on top of each bar) + algorithm's execution time.

Our OpenMP parallel implementation **Local-core(P)** decomposes *aminer* hypergraph with 27M nodes, 17M hyperedges in 91 seconds.

Application 1: Densest subgraph discovery

- A new notion of densenset sub-hypergraph.

Volume-densest subhypergraph

The **volume-densest subhypergraph** is a subhypergraph which has the largest volume-density among all subhypergraphs. The **volume-density** $\rho^N[S]$ of a subset $S \subseteq V$ of nodes in a hypergraph $\rho^N[S] = \frac{\sum_{u \in S} |N_S(u)|}{|S|}$.

- **Greedy approximation algorithm** for volume-densest subhyp. recovery is $(d_{pair}(d_{card} - 2) + 2)$ -approximate, where hyperedge-cardinality and node-pair co-occurrence (# hyperedges containing that pair) are at most d_{card} and d_{pair} , resp.

Case study: Nashville Meetup Dataset

- The degree-densest subhyp. contains casual, frequent gatherings from only one socializing group. (**Not informative**)
- The degree-densest subgraph of the clique graph captures technical events arranged by diverse, yet niche activity groups (e.g. 5 participants on avg.) (**Informative**)
- The volume-densest subhyp. captures technical events arranged by diverse and vibrant activity groups (78 participants on avg.) (**More informative**)

Application 2: Influence spreading and intervention

Initially, all nodes except one called a *seed*- are at the *susceptible* state. The seed node is initially at the *infectious* state. At each time step, each infected node infects its susceptible neighbors with probability β and then becomes *immunized*. Once a node is immunized, it is never re-infected.

- ➊ Inner-cores produced by our decomposition contain influential spreaders.
- ➋ Our decomposition produces the best order of important nodes for deleting a limited number of them while causing the maximum intervention in spreading.