



# Natural and Artificial Dynamics in GNNs



Dongqi Fu  
UIUC



Zhe Xu  
UIUC



Hanghang Tong  
UIUC



Jingrui He  
UIUC

{dongqif2, zhixu3, htong, jingrui}@illinois.edu

<https://github.com/DongqiFu/Natural-and-Artificial-Dynamics-in-GNNs-A-Tutorial>



# Contents

- **Part I – Introduction**
- **Part II – Natural Dynamics in GNNs**
- **30 mins Coffee Break**
- **Part III – Artificial Dynamics in GNNs**
- **Part IV – Challenges and Future Directions**
- **Q&A**

# Basics of graph neural networks (GNNs)

- According to [1], the general formula of GNNs can be expressed as

message-passing: information aggregation among hidden representation vectors of neighbors

$$\mathbf{a}_v^{(k)} = \text{AGGREGATE}^{(k)}(\{\mathbf{h}_u^{(k-1)} : u \in \mathcal{N}(v)\}), \quad \mathbf{h}_v^{(k)} = \text{COMBINE}^{(k)}(\mathbf{h}_v^{(k-1)}, \mathbf{a}_v^{(k)})$$

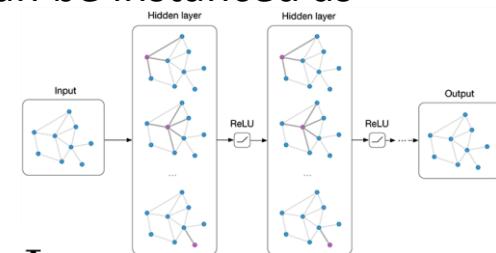
$\mathbf{h}_v^{(k)}$ : is the hidden representation of node  $v$  at the  $k$ -th layer

- For example, the graph convolutional neural network (GCN) [2] can be instanced as

$$\mathbf{h}_v^{(k)} = \text{ReLU}(\mathbf{W}^{(k-1)} \cdot \text{MEAN}\{\mathbf{h}_u^{(k-1)} : u \in \mathcal{N}(v) \cup \{v\}\})$$

with the original formula as

$$\mathbf{H}^{(k)} = \text{ReLU}(\hat{\mathbf{A}} \mathbf{H}^{(k-1)} \mathbf{W}^{(k-1)}) \quad \hat{\mathbf{A}} = \tilde{\mathbf{D}}^{-\frac{1}{2}} \tilde{\mathbf{A}} \tilde{\mathbf{D}}^{-\frac{1}{2}} \quad \tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$$



# GNNs have broad application domains



Computer Vision [1]



Natural Language Processing [2]



Recommender Systems [3]



Drug Discovery [4]

image source: <https://realpython.com/>

[1] Chen et al.: A Survey on Graph Neural Networks and Graph Transformers in Computer Vision: A Task-Oriented Perspective. CoRR 2022

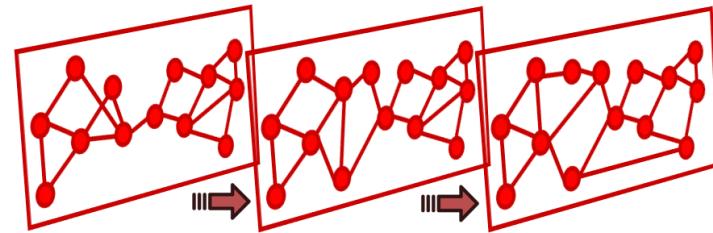
[2] Wu et al.: Graph Neural Networks for Natural Language Processing: A Survey. CoRR 2021

[3] Wang et al.: Graph Learning based Recommender Systems: A Review. IJCAI 2021

[4] Gaudelet et al.: Utilizing Graph Machine Learning within Drug Discovery and Development. Briefings in Bioinformatics 2021

# What are natural dynamics?

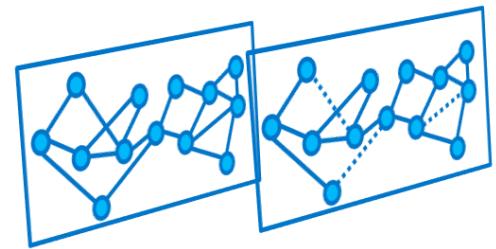
- Natural dynamics in graphs [1]
  - Input graphs have the time-evolving components, e.g.,
    - Topology Structures
    - Node-level, edge-level, and (sub)graph-level features, etc.
  - Continuous Time
    - $\mathcal{G} = \{A, e = (i, j, t, +/-)\}$
  - Discrete Time
    - $\mathcal{G} = \{A^{(1)}, A^{(2)}, \dots, A^{(T)}\}$



Evolving Graph Structures (Discrete Time Representation)

# What are artificial dynamics?

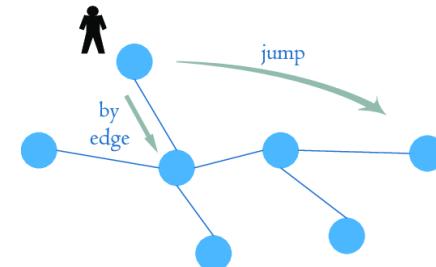
- Artificial dynamics in graphs [1], end-users
  - **change** (e.g., filter, mask, drop, or augment) **the existing** or
  - **construct the non-existing** graph-related elements, e.g.,
    - graph topology
    - node/graph attributes
    - GNN gradients, etc.
  - to realize the certain **performance upgrade**, e.g.,
    - decision accuracy
    - computation efficiency
    - model explanation, etc.



Random Edge Dropping ( $A \rightarrow \bar{A}$ )

# What are artificial dynamics?

- Artificial dynamics in graphs [1], end-users
  - **change** (e.g., filter, mask, drop, or augment) **the existing** or
  - **construct the non-existing** graph-related elements
  - to realize the certain **performance upgrade**
- In 2003, “artificial jump” [2] is proposed to adjust the graph topology for PageRank realizing the personal ranking function on graphs



# Relation between natural and artificial dynamics?

- For natural dynamics,
  - The input graph itself is a sequence of observations based on time
  - E.g., daily world wide web like Facebook, Twitter, etc.
- For artificial dynamics,
  - End-users deliberately modify the components for different interests
  - E.g., imperfect or redundant connections, missing features, etc.
- Can they be combined, i.e., **natural + artificial dynamics?**
  - Yes, when the input graph is temporal, and the modification is necessary



# How natural dynamics contribute GNNs?

- Considering natural dynamics can help graph machine learning models to capture the temporal correlations among features [1]



- Running?
- Dancing?
- Or just the static model for photography?

# How natural dynamics contribute GNNs?

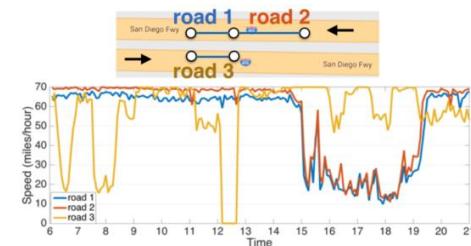
- Considering natural dynamics can help graph machine learning models to capture the temporal correlations among features [1]



- Running? 
- Dancing?
- Or just the static model for photography?

# How natural dynamics contribute GNNs?

- Considering natural dynamics can help graph machine learning models to capture the temporal correlations among features [1]
  - Motion Recognition [2]
  - Time-Series Forecasting [3]
  - Pandemic Classification [4]
  - Social Network Analysis [5]
  - Many more ...



[1] Kazemi et al.: Representation Learning for Dynamic Graphs: A Survey. JMLR (2020)

[2] Yan et al.: Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. AAAI 2018

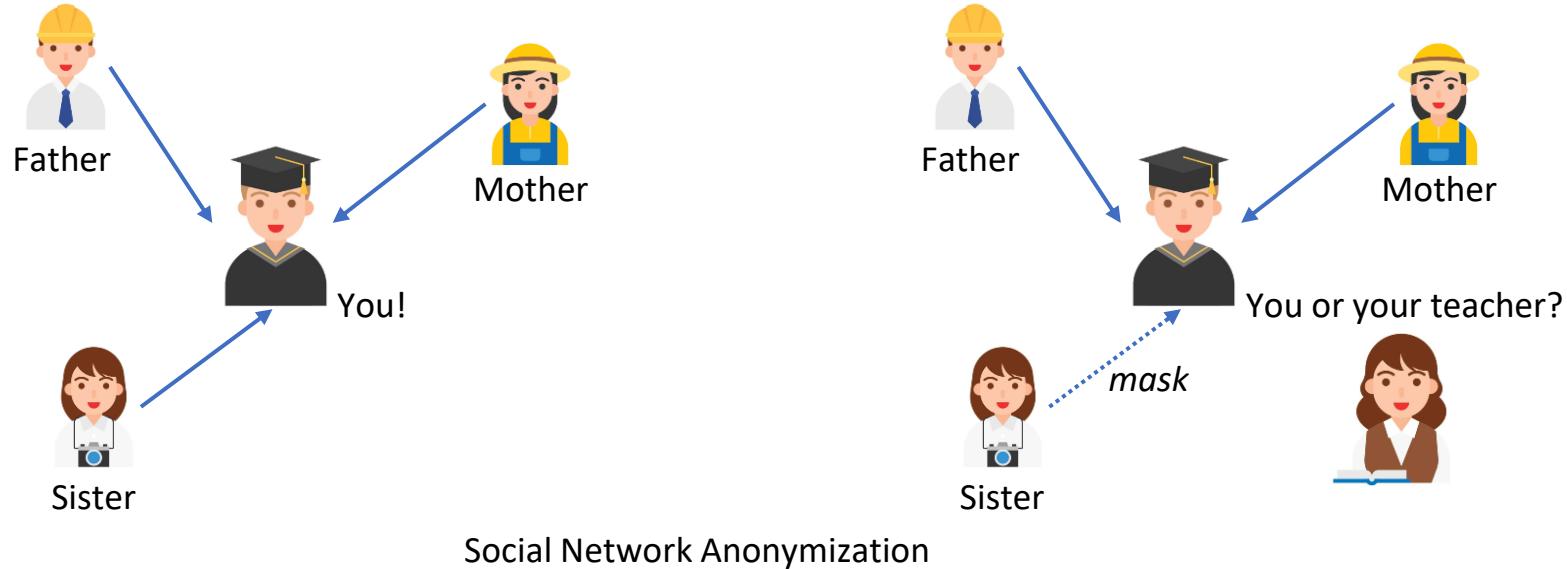
[3] Li et al.: Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting. ICLR 2018

[4] Tsiotas et al.: The Effect of Anti-COVID-19 Policies on the Evolution of the Disease: A Complex Network Analysis of the Successful Case of Greece. Physics (2020)

[5] Aggarwal et al.: Evolutionary Network Analysis: A Survey. ACM Comput. Surv. (2014)

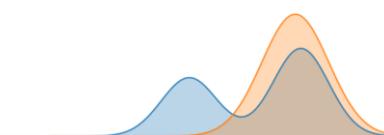
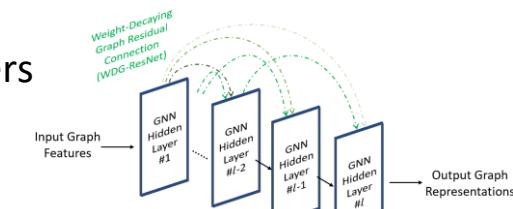
# How artificial dynamics contribute GNNs?

- Considering artificial dynamics can boost graph machine learning performance [1]



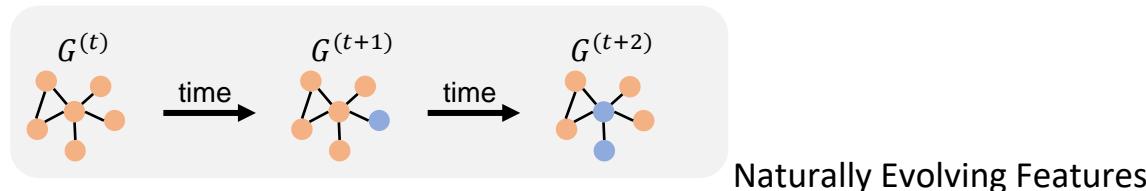
# How artificial dynamics contribute GNNs?

- Considering artificial dynamics can boost graph machine learning performance [1]
  - Privacy-Preserving [2]
    - **Permute the GNN gradients** under the differential privacy constraint
  - Decision Accuracy [3]
    - **Add dependency constraints on weight matrices** of GNN layers
  - Domain Adaption [4]
    - **Graph promoting** for large-scale pre-trained graph models on downstream tasks [4]
  - Many more ...

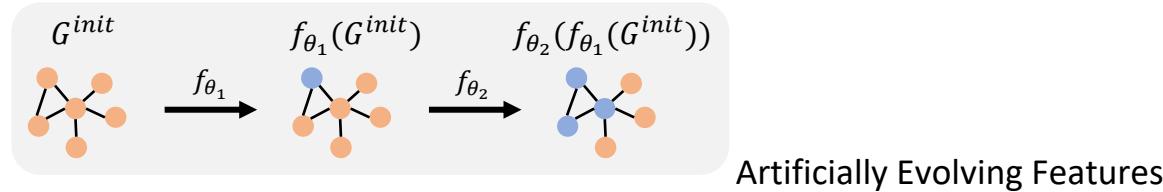


# Scope of this tutorial

- Natural dynamics in GNNs
  - We focus on **time-evolving graph structures** and node **features**

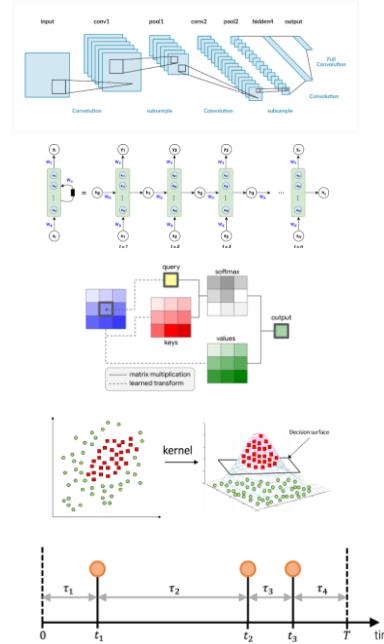


- Artificial dynamics in GNNs
  - We focus the **augmentation strategies** on graph **structures** and node **features**



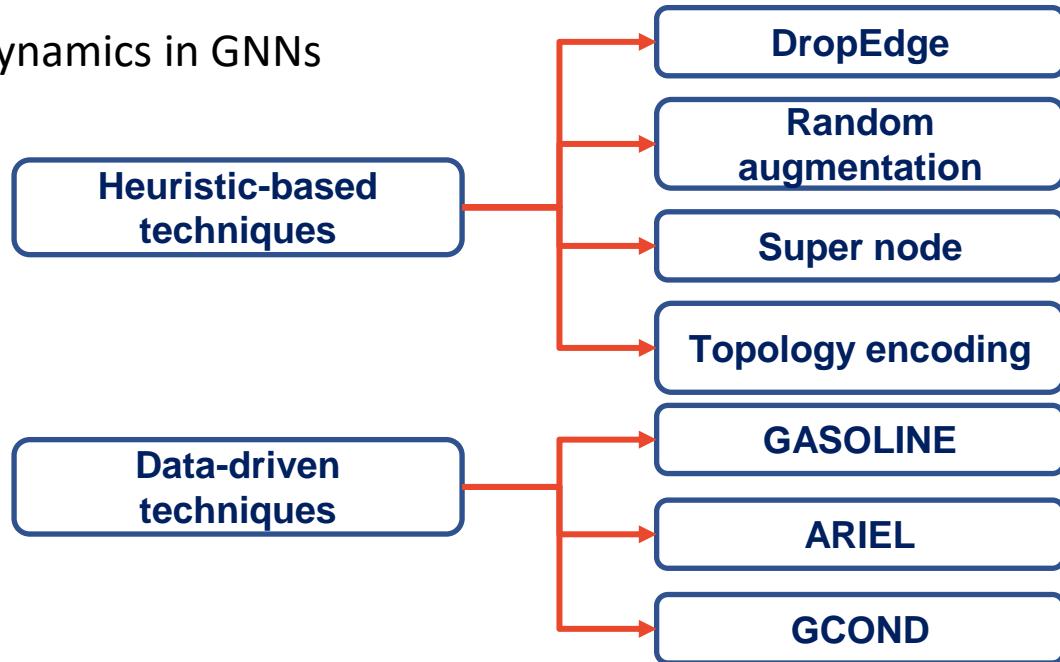
# In this tutorial

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process



# In this tutorial

- Covered works for artificial dynamics in GNNs
  - Heuristic-based Artificial Dynamics
  - Data-driven Artificial Dynamics

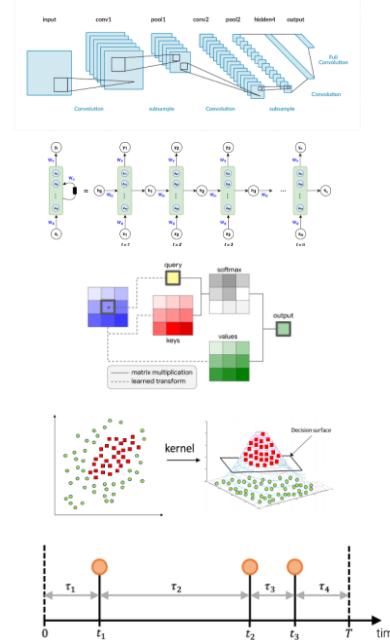


# Contents

- Part I – Introduction 
- Part II – Natural Dynamics in GNNs
- 30 mins Coffee Break
- Part III – Artificial Dynamics in GNNs
- Part IV – Challenges and Future Directions
- Q&A

# Roadmap for natural dynamics

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process

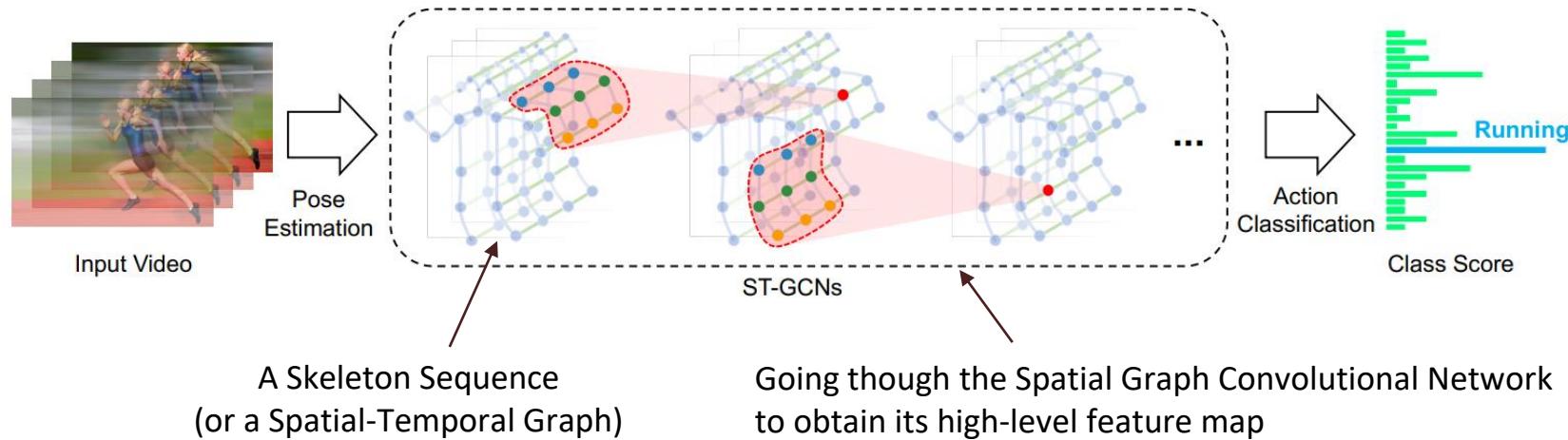


# Spatial Temporal Graph Convolutional Networks (ST-GCN) [1]

- **Task:** graph-level representation learning
- **Natural dynamic:** evolving structures and node features w.r.t time
- **Goal:** graph classification

# Spatial Temporal Graph Convolutional Networks (ST-GCN) [1]

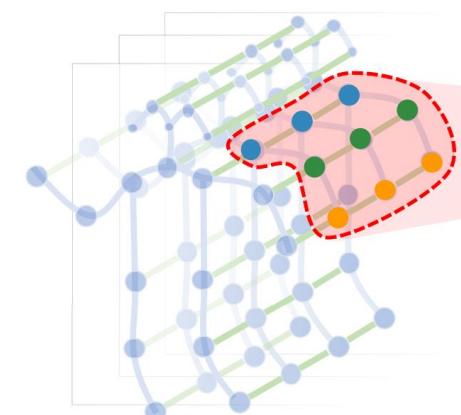
- Problem Setting
  - Skeleton-based Action Reconstruction
    - or temporal graph classification in the graph research community



[1] Sijie Yan, Yuanjun Xiong, Dahua Lin: Spatial Temporal Graph Convolutional Networks for Skeleton-Based Action Recognition. AAAI 2018

# Spatial Temporal Graph Convolutional Networks (ST-GCN) [1]

- Graph Modeling
  - $G = (V, E)$  on a skeleton sequence with  $N$  joints and  $T$  timestamps featuring both intra-body and inter-frame connections
  - $V = \{v_{ti} \mid t = 1, \dots, T, i = 1, \dots, N\}$ 
    - the  $i$ -th joint at time  $t$
  - $F(v_{ti})$ : node feature, containing coordinate vector, estimation confidence, etc.
  - $E_S = \{v_{ti}v_{tj} \mid i \neq j\}$ : human body joints
    - same  $t$
  - $E_F = \{v_{ti}v_{(t+1)i}\}$ : a particular joint  $i$ 's trajectory over time



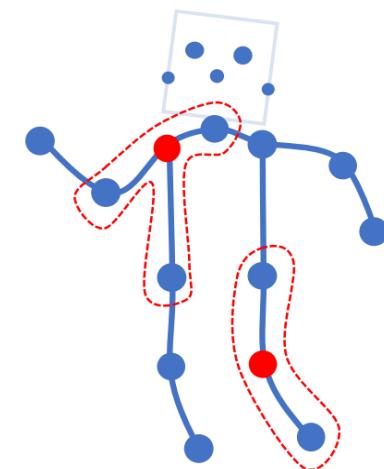
# Spatial Temporal Graph Convolutional Networks (ST-GCN) [1]

- Let's start from one single frame at  $t$ 
  - Spatial Graph Convolutional Neural Network

$$f_{out}(v_{ti}) = \sum_{v_{tj} \in B(v_{ti})} \frac{1}{Z_{ti}(v_{tj})} f_{in}(v_{tj}) \cdot \mathbf{w}(v_{ti}, v_{tj})$$

normalizing term: how many number of nodes that  
are equivalent to  $v_{tj}$ , towards  $v_{ti}$

**neighbors of  $v_{ti}$**   $B(v_{ti}) = \{v_{tj} | d(v_{tj}, v_{ti}) \leq D\}$



- which can be realized by GCN layer [2]

# Spatial Temporal Graph Convolutional Networks (ST-GCN) [1]

- Then, let's consider multiple timestamps
  - Recall the Spatial Graph Convolutional Neural Network

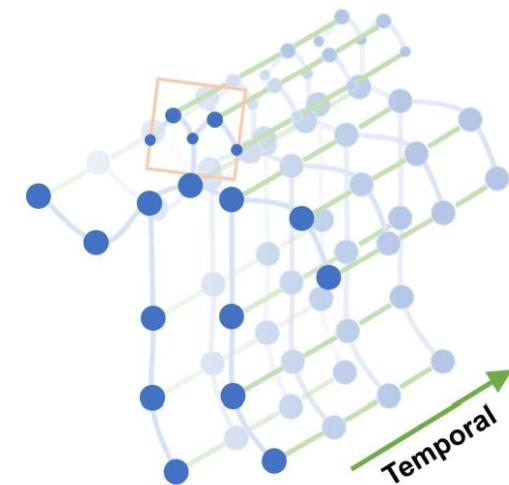
$$f_{out}(v_{ti}) = \sum_{v_{tj} \in B(v_{ti})} \frac{1}{Z_{ti}(v_{tj})} f_{in}(v_{tj}) \cdot \mathbf{w}(v_{ti}, v_{tj})$$

$$B(v_{ti}) = \{v_{tj} | d(v_{tj}, v_{ti}) \leq D\}$$

- For Spatial Temporal Graph Convolution
  - Spatial Temporal Modeling

$$B(v_{ti}) = \{v_{qj} | d(v_{tj}, v_{ti}) \leq K, |q - t| \leq \lfloor \Gamma/2 \rfloor\}$$

a hyperparameter controlling the time range



# Spatial Temporal Graph Convolutional Networks (ST-GCN) [1]

- A single frame shares the time, i.e.,

$E_S = \{v_{ti} v_{tj}\}$ : human body joints

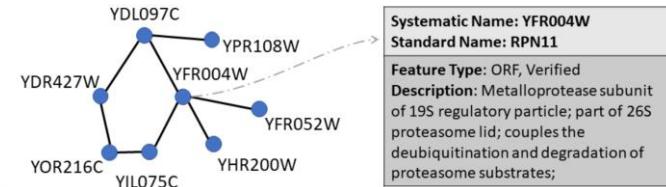
intra-snapshot edges

$$f_{out}(v_{ti}) = \sum_{v_{tj} \in B(v_{ti})} \frac{1}{Z_{ti}(v_{tj})} f_{in}(v_{tj}) \cdot \mathbf{w}(v_{ti}, v_{tj})$$

spatial graph convolution for a snapshot

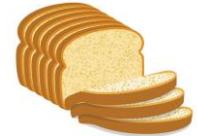
- Is that possible that they have different timestamps?

- In [2], each dynamic protein-protein interaction network has 36 continuous observations (i.e., **36 edge timestamps**)
- every 12 observations compose a metabolic cycle (i.e., **3 snapshot timestamps**), and each cycle reflects 25 mins in the real world.



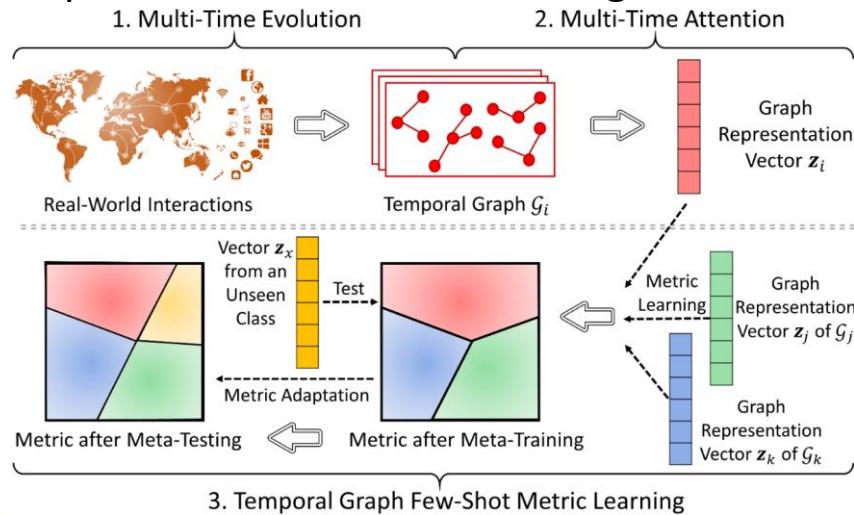
# Given multiple timestamps (e.g., edge timestamps and snapshot timestamps) in a temporal graph

- RQ1: How to integrate the multiple evolution patterns?
- RQ2: How to encode them for an embedding for temporal graph classification? What evolutions are dominating the graph similarity?
- RQ3: Labeling graph (especially temporal) is costly, how could we leverage fewer labels but effectively?



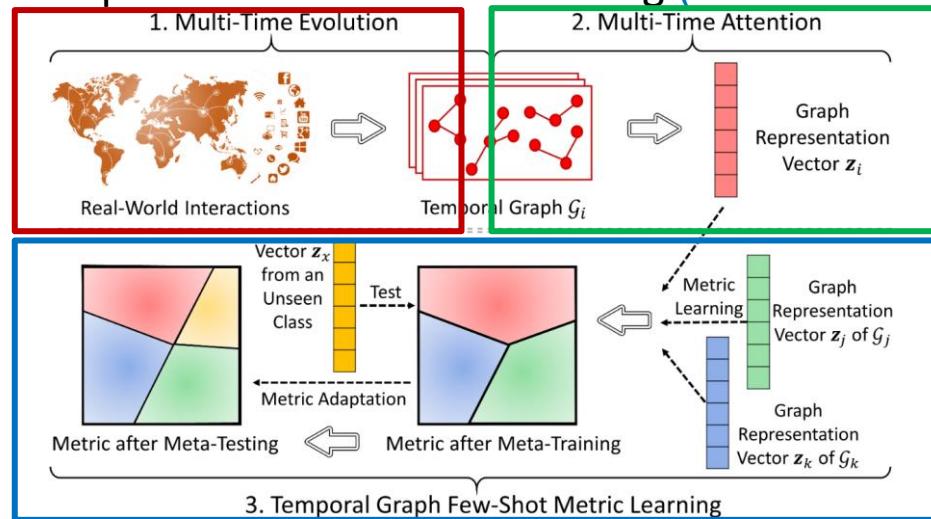
# Facing above research questions, in Temp-GFSM [1]

- Multi-Time Evolution
- Multi-Time Attention
- Temporal Graph Few-Shot Metric Learning



# In Temp-GFSM [1]

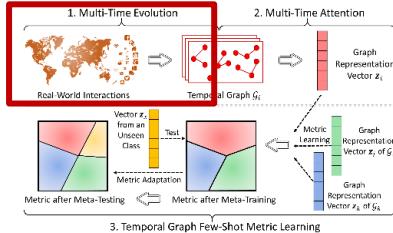
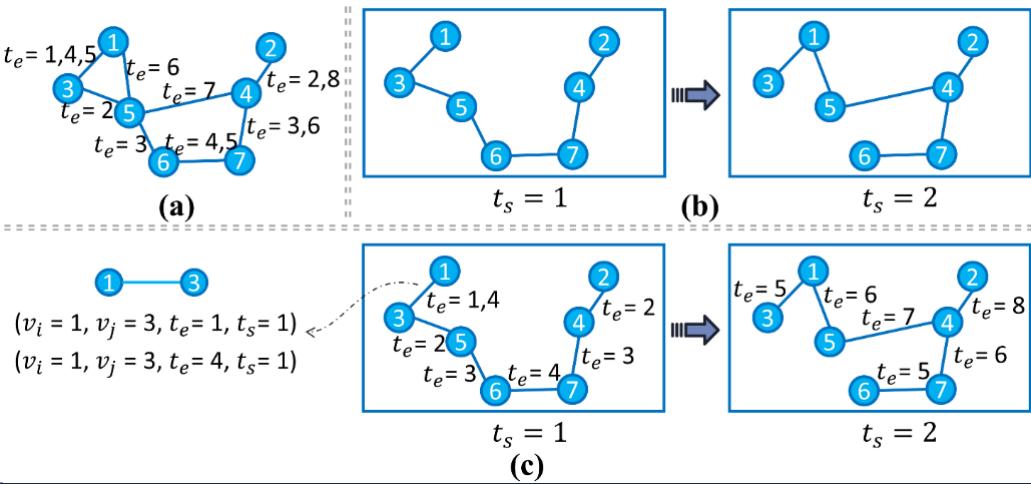
- Multi-Time Evolution (Carrying Multiple Dynamics)
- Multi-Time Attention (Weighting Multiple Dynamics)
- Temporal Graph Few-Shot Metric Learning (New Class Adaption)



[1] Dongqi Fu, Liri Fang, Ross Maciejewski, Vette I. Torvik, Jingrui He: Meta-Learned Metrics over Multi-Evolution Temporal Graphs. KDD 2022

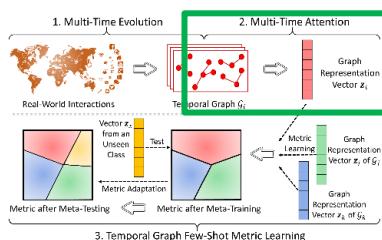
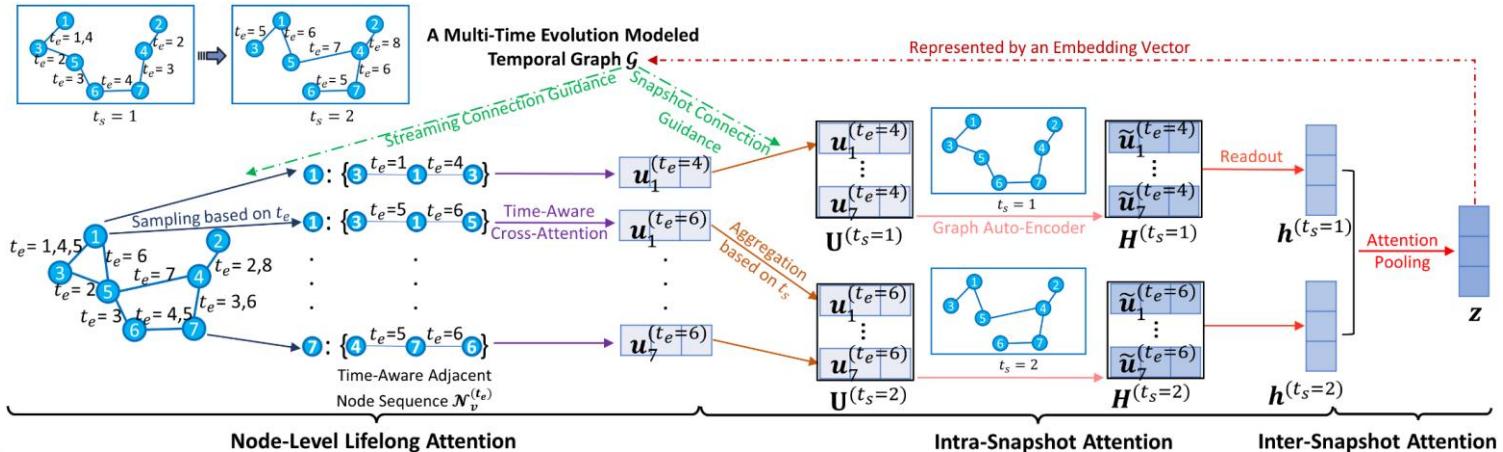
# In Multi-Time Evolution of Temp-GFMS [1]

- An edge is marked as quadruplet  $(v_i, v_j, t_e, t_s)$ , where
  - $(v_i, v_j, t_e)$  means the connection between  $v_i$  and  $v_j$  exists at time  $t_e$
  - $(v_i, v_j, t_e, t_s)$  means the event  $(v_i, v_j, t_e)$  happens in snapshot  $S^{t_s}$



# In Multi-Time Attention of Temp-GFSM [1]

- Temporal graph  $G \rightarrow$  representation vector  $Z$ 
  - Node-Level Lifelong Attention (Select Meaningful Words)
  - Intra-Snapshot Attention (Compose Supportive Sentences)
  - Inter-Snapshot Attention (Finish a Fluent Article with Paragraphs)



# Spatio-Temporal Graph Convolutional Networks (STGCN) [1]

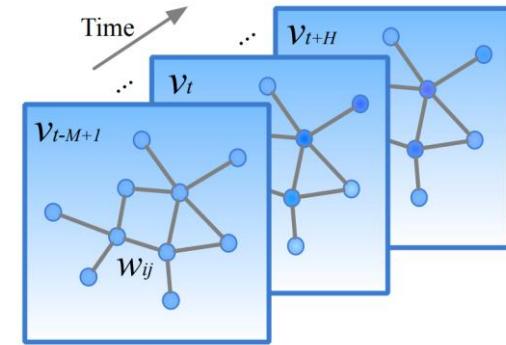
- **Task:** node-level representation learning
- **Natural dynamic:** evolving node features w.r.t time
- **Goal:** node feature prediction

# Spatio-Temporal Graph Convolutional Networks (STGCN) [1]

- A Deep Learning Framework for Traffic Forecasting
- Problem Setting
  - Traffic Flow Prediction
 

given past volumes,  
predict future volumes,  
with the latent structure

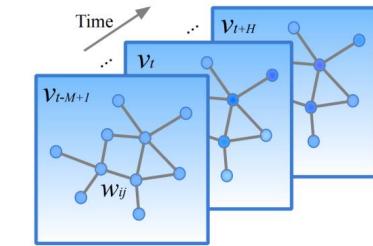
$$\arg \max \log P(v_{t+1}, \dots, v_{t+H} | v_{t-M+1}, \dots, v_t)$$
  - $v_t \in \mathbb{R}^n$ : an observation of  $n$  road segments ( $n$  nodes in graph), e.g., volume or density
  - $w \in \mathbb{R}^{n \times n}$ : adjacency matrix of road networks, the shared structure over t



# Spatio-Temporal Graph Convolutional Networks (STGCN) [1]

- Extract Spatial Features
  - Similar to previous ST-GCN [2], backbone is GCN
- Extract Temporal Features
  - Other than directly calling GCN on the catenation of temporal features,  $v^l = \{v_{t-M+1}, \dots, v_t\}$ , involve a time convolution on the time series as below

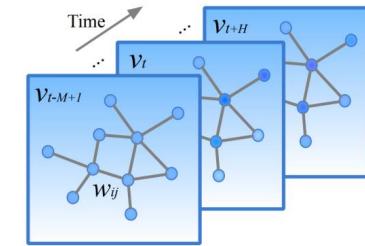
$$v^{l+1} = \Gamma_1^l *_{\mathcal{T}} \text{ReLU}(\Theta^l *_{\mathcal{G}} (\Gamma_0^l *_{\mathcal{T}} v^l))$$



$v_t \in \mathbb{R}^n$ : an observation of  $n$  road segments  
( $n$  nodes in graph), e.g., volume or density

# Spatio-Temporal Graph Convolutional Networks (STGCN) [1]

- Extract Spatial Features
  - Similar to previous ST-GCN [2], backbone is GCN
- Extract Temporal Features
  - Other than directly calling GCN on the catenation of temporal features,  $v^l = \{v_{t-M+1}, \dots, v_t\}$ , involve a time convolution on the time series as below



$$v^{l+1} = \Gamma_1^l *_{\mathcal{T}} \text{ReLU}(\Theta^l *_{\mathcal{G}} (\Gamma_0^l *_{\mathcal{T}} v^l))$$

pass a fully-connected output layer to readout  $v_{t+1}$   
 $v^l$ : stacking  $v_{t-M+1}, \dots, v_t$   
 time convolution,  
 (details next page)  
 $l$ : is the index of unit block (or layer) of STGCN, i.e.,  $v^l \rightarrow v^{l+1}$   
 graph convolution,  
 i.e., GCN

# Spatio-Temporal Graph Convolutional Networks (STGCN) [1]

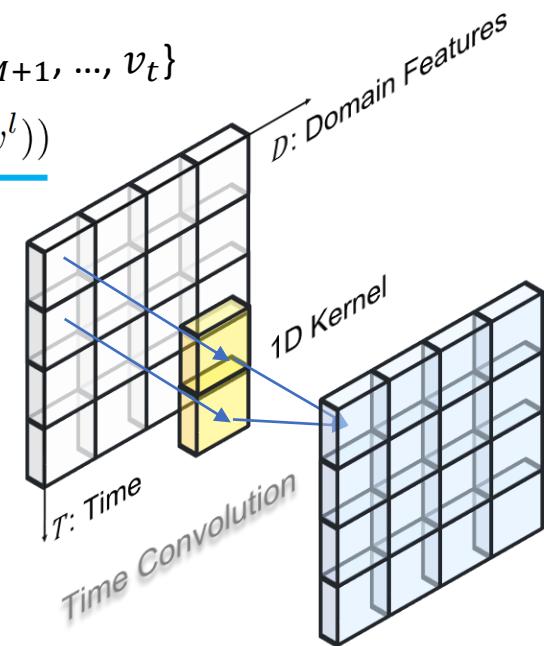
- Extract Temporal Features
  - Other than directly call GCN on the catenation  $v^l = \{v_{t-M+1}, \dots, v_t\}$

$$\text{ReLU}(\Theta^l *_{\mathcal{G}} (v^l))$$



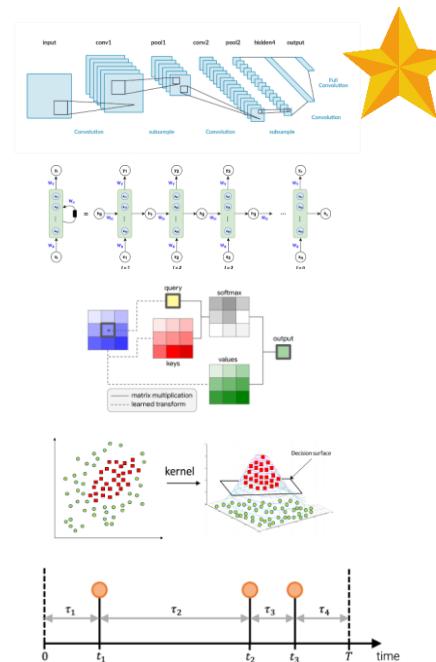
$$\text{ReLU}(\Theta^l *_{\mathcal{G}} (\Gamma_0^l *_{\mathcal{T}} v^l))$$

- A 1-D kernel along the time axis
  - Could aggregate temporal neighbors to **capture temporal behaviors** of features (e.g., traffic flows), especially for long-term time-series



# In this tutorial

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process



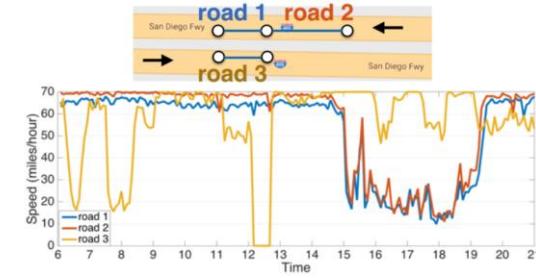
# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- **Task:** node-level representation learning
- **Natural dynamic:** evolving node features w.r.t time
- **Goal:** node feature prediction

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- Problem Definition
  - Let's still focus on Traffic Forecasting

$$[\mathbf{X}^{(t-T'+1)}, \dots, \mathbf{X}^{(t)}; \mathcal{G}] \xrightarrow{h(\cdot)} [\mathbf{X}^{(t+1)}, \dots, \mathbf{X}^{(t+T)}]$$



- But the differences from the previous discussed STGCN are:
  - What if the latent graph structure is directed?
  - How can be deal with time information other than time convolution, e.g., how to take time information recurrently?

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (1) Stationary distribution of the diffusion process

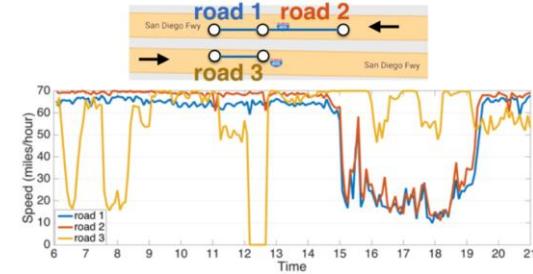
$$\mathcal{P} = \sum_{k=0}^{\infty} \alpha(1-\alpha)^k (\mathbf{D}_O^{-1} \mathbf{W})^k$$

- (2) Diffusion Convolution

$$\mathbf{X}_{:,p} \star_{\mathcal{G}} f_{\theta} = \sum_{k=0}^{K-1} \left( \theta_{k,1} (\mathbf{D}_O^{-1} \mathbf{W})^k + \theta_{k,2} (\mathbf{D}_I^{-1} \mathbf{W}^T)^k \right) \mathbf{X}_{:,p} \quad \text{for } p \in \{1, \dots, P\}$$

- (3) Diffusion Convolution Layer

$$\mathbf{H}_{:,q} = \mathbf{a} \left( \sum_{p=1}^P \mathbf{X}_{:,p} \star_{\mathcal{G}} f_{\Theta_{q,p,:,:}} \right) \quad \text{for } q \in \{1, \dots, Q\}$$



let's step into each one of those

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (1) Stationary distribution of the diffusion process can be represented as a weighted combination of infinite random walks on the graph

$$\mathcal{P} = \sum_{k=0}^{\infty} \alpha(1 - \alpha)^k (\mathbf{D}_O^{-1} \mathbf{W})^k$$

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (1) Stationary distribution of the diffusion process can be represented as a weighted combination of infinite random walks on the graph

$D_O$ : out-degree diagonal matrix

$$\mathcal{P} = \sum_{k=0}^{\infty} \alpha(1 - \alpha)^k (D_O^{-1} W)^k$$

$\mathcal{P} \in \mathbb{R}^{N \times N}$ : whose  $i$ -th row represents the likelihood of diffusion  
(i.e., personalized PageRank vector) from node  $i$

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (2) Diffusion Convolution over a graph signal  $\mathbf{X} \in \mathbb{R}^{N \times P}$  and a filter  $f_\theta$  is defined as

$$\mathbf{X}_{:,p} \star_{\mathcal{G}} f_\theta = \sum_{k=0}^{K-1} \left( \theta_{k,1} (\mathbf{D}_O^{-1} \mathbf{W})^k + \theta_{k,2} (\mathbf{D}_I^{-1} \mathbf{W}^\top)^k \right) \mathbf{X}_{:,p} \quad \text{for } p \in \{1, \dots, P\}$$

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (2) Diffusion Convolution **over** a graph signal  $\mathbf{X} \in \mathbb{R}^{N \times P}$  and a filter  $f_\theta$  is defined as

feature matrix,  $N$  is the num of node,  
 $P$  is the node feature dimension

learnable parameter,  
i.e., weight matrices

$$\underline{\mathbf{X}_{:,p} \star_{\mathcal{G}} f_\theta} = \sum_{k=0}^{K-1} \left( \theta_{k,1} (\mathbf{D}_O^{-1} \mathbf{W})^k + \theta_{k,2} (\mathbf{D}_I^{-1} \mathbf{W}^\top)^k \right) \mathbf{X}_{:,p} \quad \text{for } p \in \{1, \dots, P\}$$

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (2) Diffusion Convolution over a graph signal  $X \in \mathbb{R}^{N \times P}$  and a filter  $f_\theta$  is defined as

feature matrix,  $N$  is the num of node,  
 $P$  is the node feature dimension

learnable parameter,  
i.e., weight matrices

$$X_{:,p} \star_{\mathcal{G}} f_\theta = \sum_{k=0}^{K-1} \left( \underbrace{\theta_{k,1}}_{\text{two sets of parameters from } f_\theta} \underbrace{\left( D_O^{-1} W \right)^k}_{\text{out-degree based diffusion}} + \underbrace{\theta_{k,2}}_{\text{two sets of parameters from } f_\theta} \underbrace{\left( D_I^{-1} W^\top \right)^k}_{\text{in-degree based diffusion}} \right) \underbrace{X_{:,p}}_{\text{feature matrix}} \quad \text{for } p \in \{1, \dots, P\}$$

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (3) Diffusion Convolution Layer that maps  $P$ -dimensional features to  $Q$ -dimensional outputs

$$\mathbf{H}_{:,q} = \mathbf{a} \left( \sum_{p=1}^P \mathbf{X}_{:,p} \star_{\mathcal{G}} f_{\Theta_{q,p,:,:}} \right) \quad \text{for } q \in \{1, \dots, Q\}$$

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- For directed graph structure
  - (3) Diffusion Convolution Layer that maps  $P$ -dimensional features to  $Q$ -dimensional outputs

$X \in \mathbb{R}^{N \times P}$ : input

$H \in \mathbb{R}^{N \times Q}$ : output

$$H_{:,q} = a \left( \sum_{p=1}^P X_{:,p} \star_{\mathcal{G}} f_{\Theta_{q,p},:,:} \right)$$

for  $q \in \{1, \dots, Q\}$

activation function (e.g., ReLU, Sigmoid)

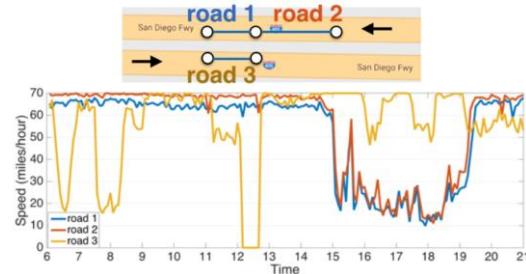
indexing for in-degree or out-degree diffusion

indexing for  $k$

$$X_{:,p} \star_{\mathcal{G}} f_{\theta} = \sum_{k=0}^{K-1} \left( \theta_{k,1} (D_O^{-1} W)^k + \theta_{k,2} (D_I^{-1} W^\top)^k \right) X_{:,p} \quad \text{for } p \in \{1, \dots, P\}$$

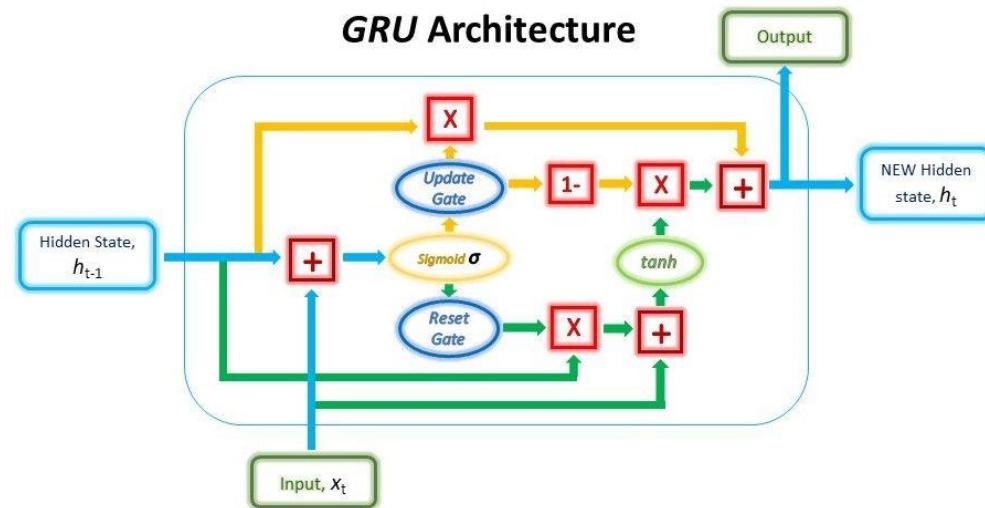
# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- After for directed graph structure,
  - (1) Stationary distribution of the diffusion process
  - (2) Diffusion Convolution
  - (3) Diffusion Convolution Layer
- How to take time information recurrently?
  - Temporal Dynamics Modeling [1]



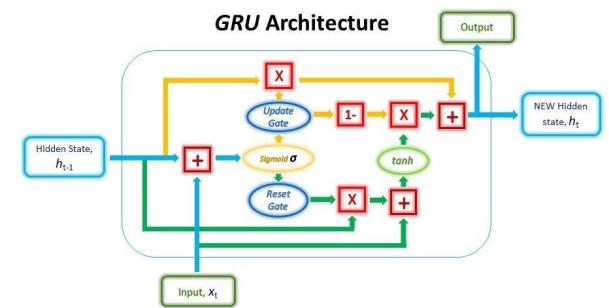
# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- Temporal Dynamics Modeling
  - Adopt the logic from Gated Recurrent Units (GRU)[2], i.e., make GRU take structured information
  - Input
  - Reset Gate
  - Update Gate
  - Output



# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- Temporal Dynamics Modeling
  - Adopt the logic from Gated Recurrent Units (GRU)[2], i.e., make GRU take structured information
  - Input  $\mathbf{X}^{(t)}, \mathbf{H}^{(t-1)}$
  - Reset Gate  $\mathbf{r}^{(t)} = \sigma(\Theta_r \star_{\mathcal{G}} [\mathbf{X}^{(t)}, \mathbf{H}^{(t-1)}] + \mathbf{b}_r)$
  - Update Gate  $\mathbf{u}^{(t)} = \sigma(\Theta_u \star_{\mathcal{G}} [\mathbf{X}^{(t)}, \mathbf{H}^{(t-1)}] + \mathbf{b}_u)$
  - Output  $\mathbf{H}^{(t)} = \mathbf{u}^{(t)} \odot \mathbf{H}^{(t-1)} + (1 - \mathbf{u}^{(t)}) \odot \mathbf{C}^{(t)}$   
 $\mathbf{C}^{(t)} = \tanh(\Theta_C \star_{\mathcal{G}} [\mathbf{X}^{(t)}, (\mathbf{r}^{(t)} \odot \mathbf{H}^{(t-1)})] + \mathbf{b}_c)$

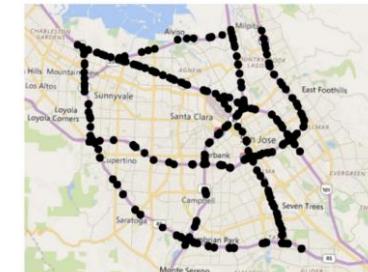
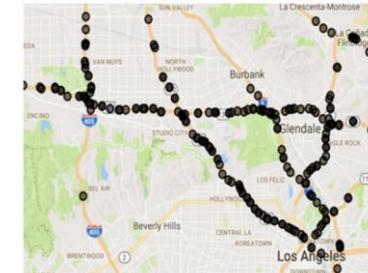


diffusion convolution w. different

weight parameters

# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- Experiments
  - Datasets
    - METR-LA: Highway of Los Angeles
      - 207 sensors (traffic speed)
      - May 1<sup>st</sup> 2012 to Jun 30<sup>th</sup> 2012
    - PEMS-BAY: Highway in Bay Area of California
      - 325 sensors (traffic speed)
      - Jan 1<sup>st</sup> 2017 to May 31th 2017

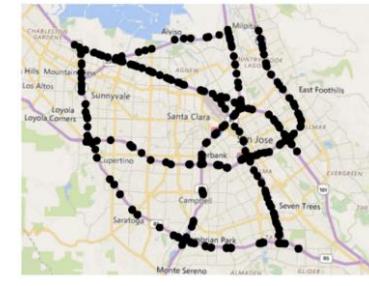
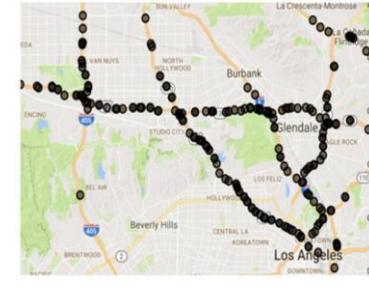


# Diffusion Convolutional Recurrent Neural Network (DCRNN) [1]

- Experiments

- Performance (traffic speed forecasting)

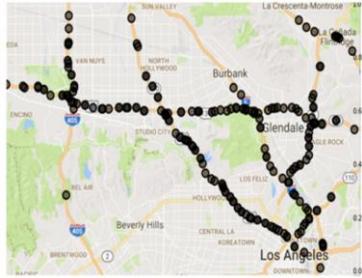
	$T$	Metric	HA	ARIMA <sub>Kal</sub>	VAR	SVR	FNN	FC-LSTM	DCRNN
METR-LA	15 min	MAE	4.16	3.99	4.42	3.99	3.99	3.44	<b>2.77</b>
		RMSE	7.80	8.21	7.89	8.45	7.94	6.30	<b>5.38</b>
		MAPE	13.0%	9.6%	10.2%	9.3%	9.9%	9.6%	<b>7.3%</b>
	30 min	MAE	4.16	5.15	5.41	5.05	4.23	3.77	<b>3.15</b>
		RMSE	7.80	10.45	9.13	10.87	8.17	7.23	<b>6.45</b>
		MAPE	13.0%	12.7%	12.7%	12.1%	12.9%	10.9%	<b>8.8%</b>
	1 hour	MAE	4.16	6.90	6.52	6.72	4.49	4.37	<b>3.60</b>
		RMSE	7.80	13.23	10.11	13.76	8.69	8.69	<b>7.59</b>
		MAPE	13.0%	17.4%	15.8%	16.7%	14.0%	13.2%	<b>10.5%</b>
PEMS-BAY	15 min	MAE	2.88	1.62	1.74	1.85	2.20	2.05	<b>1.38</b>
		RMSE	5.59	3.30	3.16	3.59	4.42	4.19	<b>2.95</b>
		MAPE	6.8%	3.5%	3.6%	3.8%	5.19%	4.8%	<b>2.9%</b>
	30 min	MAE	2.88	2.33	2.32	2.48	2.30	2.20	<b>1.74</b>
		RMSE	5.59	4.76	4.25	5.18	4.63	4.55	<b>3.97</b>
		MAPE	6.8%	5.4%	5.0%	5.5%	5.43%	5.2%	<b>3.9%</b>
	1 hour	MAE	2.88	3.38	2.93	3.28	2.46	2.37	<b>2.07</b>
		RMSE	5.59	6.50	5.44	7.08	4.98	4.96	<b>4.74</b>
		MAPE	6.8%	8.3%	6.5%	8.0%	5.89%	5.7%	<b>4.9%</b>



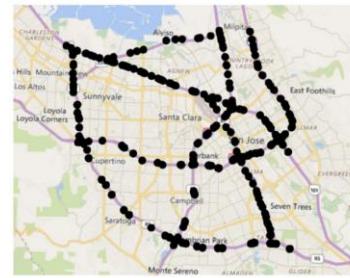
# The latent graph structure may be not easy to observe

- In DCRNN[1], the adjacency is hand-crafted

$$W_{ij} = \exp\left(-\frac{\text{dist}(v_i, v_j)^2}{\sigma^2}\right)$$



(a) METR-LA



(b) PEMS-BAY

- Could we find another way to extract that latent structure?

# Could we find another way to extract that latent structure?

- Discrete Graph Structure Learning for Forecasting Multiple Time Series (GTS) [2]
  - Focusing on the same problem (i.e., traffic forecasting) and the same diffusion convolution structure as DCRNN [1]

$$[\mathbf{X}^{(t-T'+1)}, \dots, \mathbf{X}^{(t)}; \mathcal{G}] \xrightarrow{h(\cdot)} [\mathbf{X}^{(t+1)}, \dots, \mathbf{X}^{(t+T)}]$$

- But set the adjacency matrix as a variable to learn

$$A_{ij} = \text{sigmoid}\left(\left(\log\left(\theta_{ij}/(1 - \theta_{ij})\right) + (g_{ij}^1 - g_{ij}^2)\right)/s\right)$$

temperature

$$\theta_{ij} = \text{FC}(\text{FC}(z^i \| z^j))$$

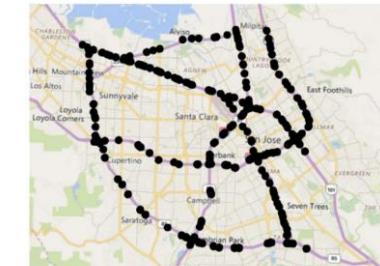
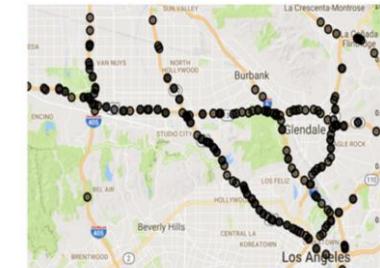
samples from a given Gumbel distribution

$$z^i = \text{FC}(\text{vec}(\text{Conv}(X^i))) \quad X^i: \text{the } i\text{-th node over all features and timestamps}$$

# In the same datasets with [1], i.e., METR-LA and PEMs-BAY

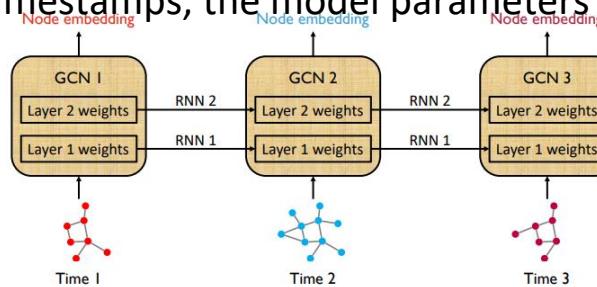
- Experiments
  - Performance (traffic speed forecasting)

	Metric	FNN	LSTM	DCRNN	LDS	NRI	GTSv	GTS
15 min	MAE ( $\times 10^{-3}$ )	1.23	1.02	0.71	0.49	0.66	0.26	<b>0.24</b>
	RMSE ( $\times 10^{-2}$ )	1.28	1.63	1.42	1.26	0.27	0.20	<b>0.19</b>
	MAPE	0.20%	0.21%	0.09%	0.07%	0.14%	0.05%	<b>0.04%</b>
30 min	MAE ( $\times 10^{-3}$ )	1.42	1.11	1.08	0.81	0.71	0.31	<b>0.30</b>
	RMSE ( $\times 10^{-2}$ )	1.81	2.06	1.91	1.79	0.30	0.23	<b>0.22</b>
	MAPE	0.23%	0.20%	0.15%	0.12%	0.15%	<b>0.05%</b>	<b>0.05%</b>
60 min	MAE ( $\times 10^{-3}$ )	1.88	1.79	1.78	1.45	0.83	<b>0.39</b>	0.41
	RMSE ( $\times 10^{-2}$ )	2.58	2.75	2.65	2.54	0.46	0.32	<b>0.30</b>
	MAPE	0.29%	0.27%	0.24%	0.22%	0.17%	<b>0.07%</b>	<b>0.07%</b>



# In DCRNN [1] or GTS [2],

- The adjacency is fixed with evolving node features
- What if the adjacency is also evolving?
- For evolving structures and features
  - EvolveGCN [3] is proposed to adapt model parameters, i.e.,
    - Each time has its own GCN model with its  $A^{(t)}$  and  $H^{(t)}$
    - Cross timestamps, the model parameters are dependent, e.g.,  $W_t^{(l)} = LSTM(W_{t-1}^{(l)})$



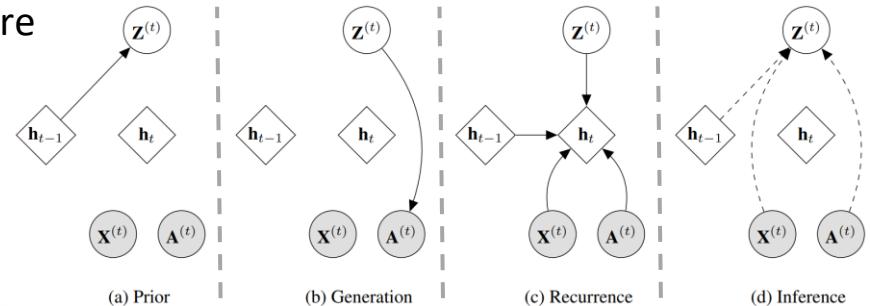
[1] Yaguang Li, Rose Yu, Cyrus Shahabi, Yan Liu: Diffusion Convolutional Recurrent Neural Network: Data-Driven Traffic Forecasting. ICLR 2018

[2] Chao Shang, Jie Chen, Jinbo Bi: Discrete Graph Structure Learning for Forecasting Multiple Time Series. ICLR 2021

[3] Pareja et al.: EvolveGCN: Evolving Graph Convolutional Networks for Dynamic Graphs. AAAI 2020

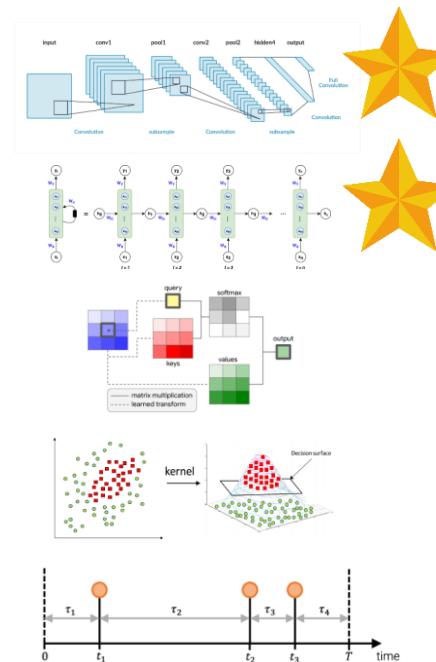
# In DCRNN [1] or GTS [2],

- The adjacency is fixed with evolving node features
- ~~What if the adjacency is also evolving?~~
- Can we also predict the future adjacency?
- For evolving structures and features
  - VGRNN [3] is proposed to learn the variational posterior distribution of evolving adjacency structures together, in the RNN structure



# In this tutorial

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process

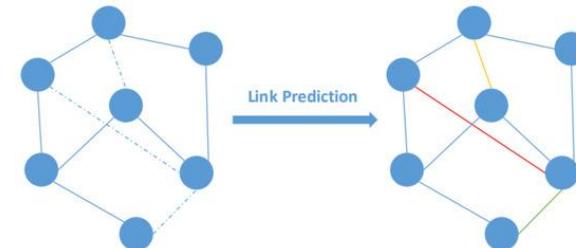


# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- **Task:** node-level representation learning
- **Natural dynamic:** evolving graph structures and node features w.r.t time
- **Goal:** link prediction

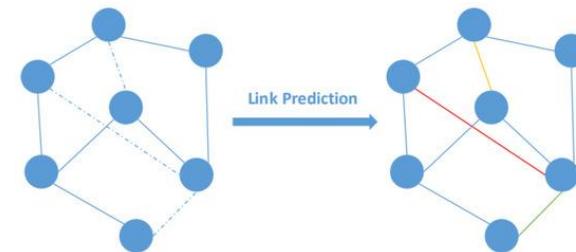
# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Problem Definition (Link Prediction)
  - Given a series of graph snapshots  $\{\mathcal{G}^1, \dots, \mathcal{G}^T\}$ , and  $\mathcal{G}^t = (\mathbf{A}^t, \mathbf{X}^t)$ , DySAT [1] aims to learn the node representation  $e_v^t$  for each node  $v$  at timestamps  $t = \{1, 2, \dots, T\}$
  - Then, the latest time  $e_v^T$  and  $e_u^T$  are used to decide if there is an edge links node  $v$  and node  $u$  at time  $T+1$



# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Structural Self-Attention
  - Apply attention in one single timestamp
  - Topological neighbors
- Temporal Self-Attention
  - Apply attention across timestamps
  - Temporal neighbors



# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Structural Self-Attention
  - At a single timestamp  $t$ , the superscript of  $t$  is omitted in the following equation

$z_v$ : hidden representation of node  $v$   
after structural attention

$$z_v = \sigma \left( \sum_{u \in N_v} \alpha_{uv} W^s x_u \right), \quad \alpha_{uv} = \frac{\exp(e_{uv})}{\sum_{w \in N_v} \exp(e_{wv})}$$

neighbors in  $A^t$

$A_{uv}$ : adjacency

$$e_{uv} = \sigma \left( A_{uv} \cdot a^T [W^s x_u || W^s x_v] \right) \quad \forall (u, v) \in \mathcal{E}$$

$a$ : learnable weight matrix (e.g., MLP)

$W^s$ : projection weight matrix for structural attention

# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Temporal Self-Attention
  - Who are temporal neighbors for a certain node?
  - Temporal neighbors for node  $v$  consist of its **historical behaviors**

# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Temporal Self-Attention
  - Temporal neighbors for node  $v$  consist of its historical behaviors

Node  $v$ 's time embedding,  $\mathbb{R}^{T \times F'}$

$$Z_v = \beta_v(X_v W_v), \quad e_v^{ij} = \left( \frac{((X_v W_q)(X_v W_k)^T)_{ij}}{\sqrt{F'}} + M_{ij} \right)$$

$X_v$ : concatenation of  $x_v^1, x_v^2, \dots, x_v^T$

$\mathbb{R}^{T \times T}, \mathbb{R}^{T \times D'}, \mathbb{R}^{D' \times F'}$

$\beta_v^{ij} = \frac{\exp(e_v^{ij})}{\sum_{k=1}^T \exp(e_v^{ik})}$ , attention weights,  $i$  and  $j$  are two timestamps

$F'$ : normalization factor

A matrix enforcing auto-regressive manner, details next page

# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Temporal Self-Attention
  - Temporal neighbors for node  $v$  are its historical behaviors

$$Z_v = \beta_v(X_v W_v), \quad \beta_v^{ij} = \frac{\exp(e_v^{ij})}{\sum_{k=1}^T \exp(e_v^{ik})},$$

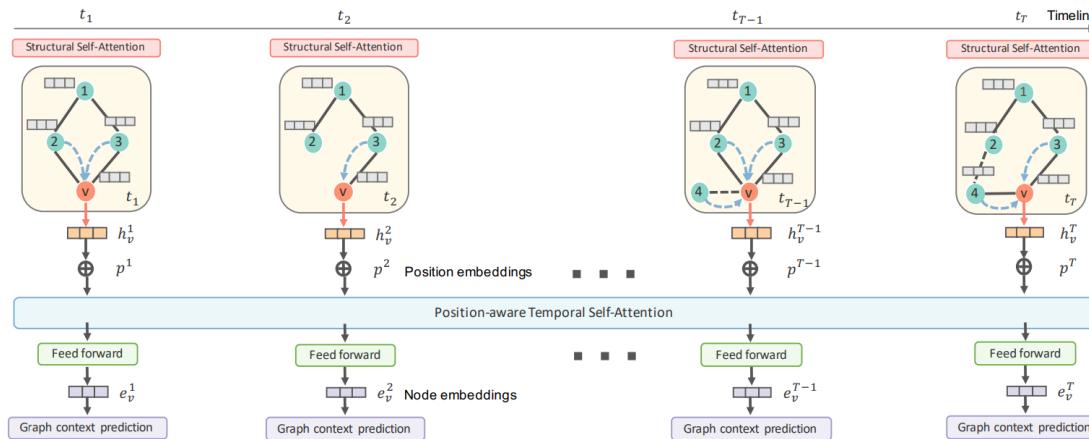
$$e_v^{ij} = \left( \frac{((X_v W_q)(X_v W_k)^T)_{ij}}{\sqrt{F'}} + M_{ij} \right)$$

$$M_{ij} = \begin{cases} 0, & i \leq j \\ -\infty, & \text{otherwise} \end{cases}$$

when  $M_{ij} = -\infty$ ,  $\beta_v^{ij} = 0$ , which switches off the attention from timestamp  $i$  to  $j$

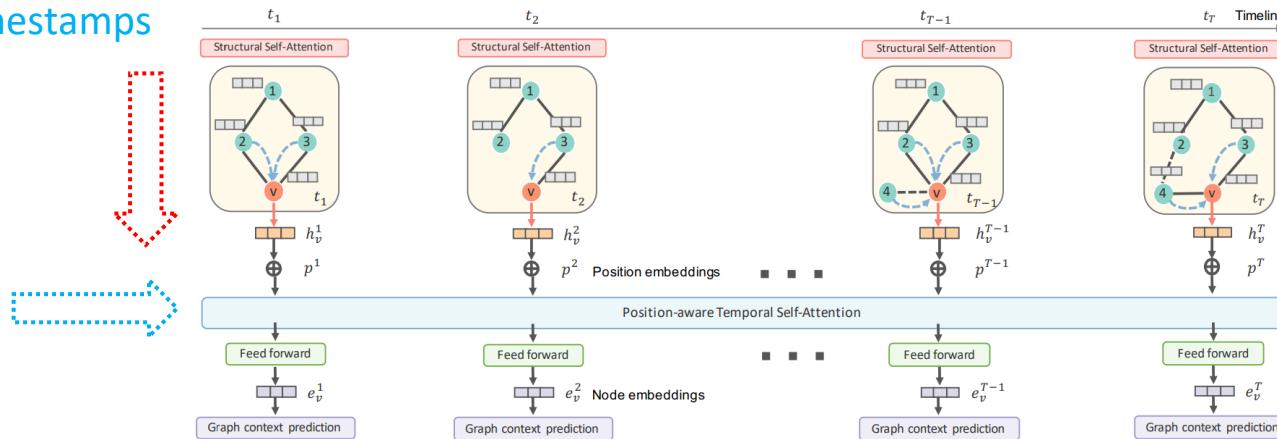
# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Structural + Temporal Self-Attention
  - Obtain structural encoding independently at each timestamp  $t$
  - Then, temporal self-attention take the structural encoding as input to attend over timestamps



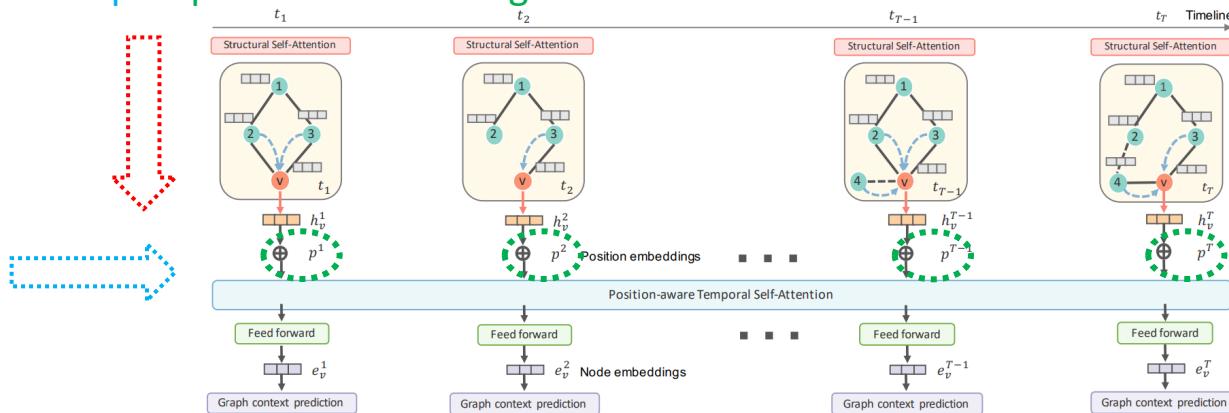
# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Structural + Temporal Self-Attention
  - Obtain structural encoding independently at each timestamp  $t$
  - Then, temporal self-attention take the structural encoding as input to attend over timestamps



# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

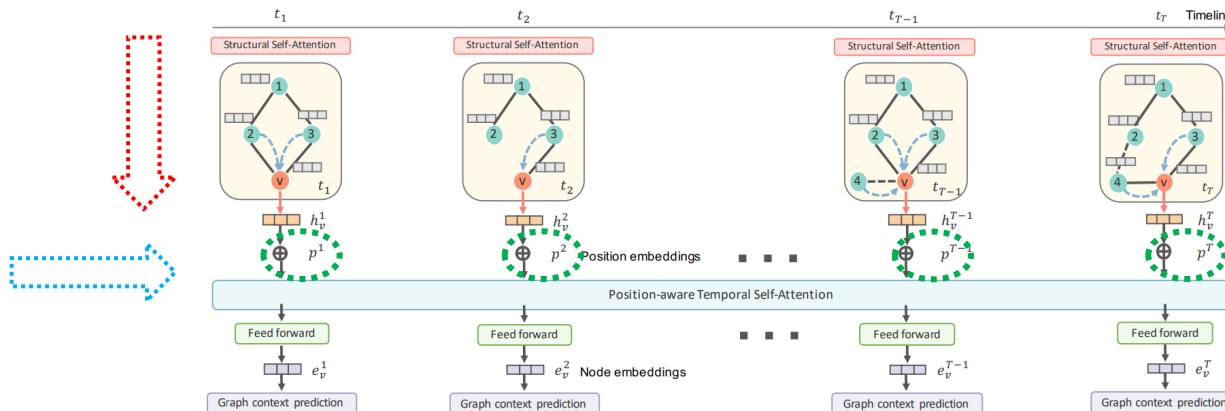
- Structural + Temporal Self-Attention
  - Obtain structural encoding independently at each timestamp  $t$
  - Then, temporal self-attention take the structural encoding as input to attend over timestamps + positional encoding



# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

- Structural + Temporal Self-Attention
  - Add the absolute temporal position of each snapshot

$$h_v^1 + p^1, h_v^2 + p^2, \dots, h_v^T + p^T$$



[1] Aravind Sankar, Yanhong Wu, Liang Gou, Wei Zhang, Hao Yang: DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention Networks.  
WSDM 2020

# DySAT: Deep Neural Representation Learning on Dynamic Graphs via Self-Attention [1]

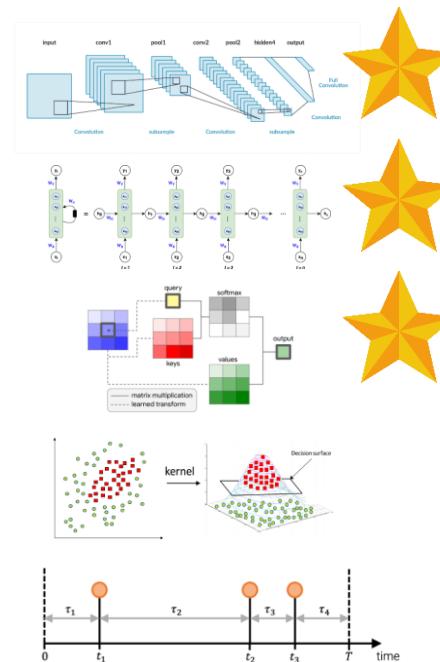
- Positional Encoding

I	→	0	P <sub>00</sub> =sin(0) = 0	P <sub>01</sub> =cos(0) = 1	P <sub>02</sub> =sin(0) = 0	P <sub>03</sub> =cos(0) = 1
am	→	1	P <sub>10</sub> =sin(1/1) = 0.84	P <sub>11</sub> =cos(1/1) = 0.54	P <sub>12</sub> =sin(1/10) = 0.10	P <sub>13</sub> =cos(1/10) = 1.0
a	→	2	P <sub>20</sub> =sin(2/1) = 0.91	P <sub>21</sub> =cos(2/1) = -0.42	P <sub>22</sub> =sin(2/10) = 0.20	P <sub>23</sub> =cos(2/10) = 0.98
Robot	→	3	P <sub>30</sub> =sin(3/1) = 0.14	P <sub>31</sub> =cos(3/1) = -0.99	P <sub>32</sub> =sin(3/10) = 0.30	P <sub>33</sub> =cos(3/10) = 0.96

image source: <https://machinelearningmastery.com/a-gentle-introduction-to-positional-encoding-in-transformer-models-part-1/>

# In this tutorial

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process



# Back to the positional encoding

I	→	0	$P_{00}=\sin(0) = 0$	$P_{01}=\cos(0) = 1$	$P_{02}=\sin(0) = 0$	$P_{03}=\cos(0) = 1$
am	→	1	$P_{10}=\sin(1/1) = 0.84$	$P_{11}=\cos(1/1) = 0.54$	$P_{12}=\sin(1/10) = 0.10$	$P_{13}=\cos(1/10) = 1.0$
a	→	2	$P_{20}=\sin(2/1) = 0.91$	$P_{21}=\cos(2/1) = -0.42$	$P_{22}=\sin(2/10) = 0.20$	$P_{23}=\cos(2/10) = 0.98$
Robot	→	3	$P_{30}=\sin(3/1) = 0.14$	$P_{31}=\cos(3/1) = -0.99$	$P_{32}=\sin(3/10) = 0.30$	$P_{33}=\cos(3/10) = 0.96$

- Do we have other options?
  - A concurrent method [1] with DySAT [2] proposes the time kernel function to record the time features

# Inductive Representation Learning on Temporal Graphs (TGAT) [1]

- **Task:** node-level representation learning
- **Natural dynamic:** evolving graph structures and node features w.r.t time
- **Goal:** link prediction, node classification

# Inductive Representation Learning on Temporal Graphs (TGAT) [1]

- In [1], **kernel function** is proposed to map time  $t$  to a finite dimensional representation vector

$$t \mapsto \Phi_d(t) := \sqrt{\frac{1}{d}} [\cos(\omega_1 t), \sin(\omega_1 t), \dots, \cos(\omega_d t), \sin(\omega_d t)]$$

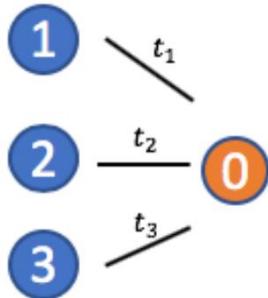
$\omega_1, \omega_2, \dots, \omega_d$  are learnable parameters

# Inductive Representation Learning on Temporal Graphs (TGAT) [1]

- Suppose there is a target node  $v_0$  at time  $t$ , which needs to attend over its spatial-temporal neighbors

$$\mathcal{N}(v_0; t) = \{v_1, \dots, v_N\}$$

- For each node  $v_i$ ,  $v_0$  connects with it previously at a time  $t_i$ , i.e.,  $t_i < t$



[1] Da Xu, Chuanwei Ruan, Evren Körpeoglu, Sushant Kumar, Kannan Achan: Inductive representation learning on temporal graphs. ICLR 2020

# Inductive Representation Learning on Temporal Graphs (TGAT) [1]

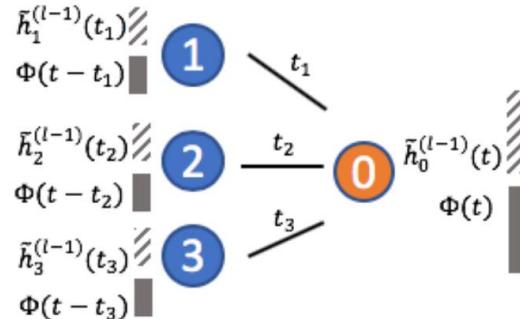
- With  $\mathcal{N}(v_0; t) = \{v_1, \dots, v_N\}$ ,  $t_i < t$
- First, append the position encoding by time kernel functions, to form  $\mathbf{Z}(t)$

$$\mathbf{Z}(t) = [\tilde{\mathbf{h}}_0^{(l-1)}(t) \parallel \Phi_{d_T}(0), \tilde{\mathbf{h}}_1^{(l-1)}(t_1) \parallel \Phi_{d_T}(t - t_1), \dots, \tilde{\mathbf{h}}_N^{(l-1)}(t_N) \parallel \Phi_{d_T}(t - t_N)]^\top$$

$\mathbf{Z}(t)$  is the intermediate step of getting the node embedding  $\tilde{\mathbf{h}}_0^{(l)}$  of  $v_0$  at the  $l$ -th layer of TGAT

node embedding of  $v_0$  at the  $(l-1)$ -th layer of TGAT, the 0-th layer is the initial input feature

time kernel function



# Inductive Representation Learning on Temporal Graphs (TGAT) [1]

- With  $\mathcal{N}(v_0; t) = \{v_1, \dots, v_N\}$ ,  $t_i < t$
- Append the position encoding by time kernel functions

$$\mathbf{Z}(t) = \left[ \tilde{\mathbf{h}}_0^{(l-1)}(t) || \Phi_{d_T}(0), \tilde{\mathbf{h}}_1^{(l-1)}(t_1) || \Phi_{d_T}(t - t_1), \dots, \tilde{\mathbf{h}}_N^{(l-1)}(t_N) || \Phi_{d_T}(t - t_N) \right]^\top$$

- Self-Attention**       $\mathbf{q}(t) = [\mathbf{Z}(t)]_0 \mathbf{W}_Q$   
 $\mathbf{K}(t) = [\mathbf{Z}(t)]_{1:N} \mathbf{W}_K$        $\mathbf{h}(t) = \text{Attn}(\mathbf{q}(t), \mathbf{K}(t), \mathbf{V}(t))$   
 $\mathbf{V}(t) = [\mathbf{Z}(t)]_{1:N} \mathbf{W}_V$
- Readout**       $\tilde{\mathbf{h}}_0^{(l)}(t) = \text{FFN}(\mathbf{h}(t) || \mathbf{x}_0) \equiv \text{ReLU}([\mathbf{h}(t) || \mathbf{x}_0] \mathbf{W}_0^{(l)} + \mathbf{b}_0^{(l)}) \mathbf{W}_1^{(l)} + \mathbf{b}_1^{(l)}$

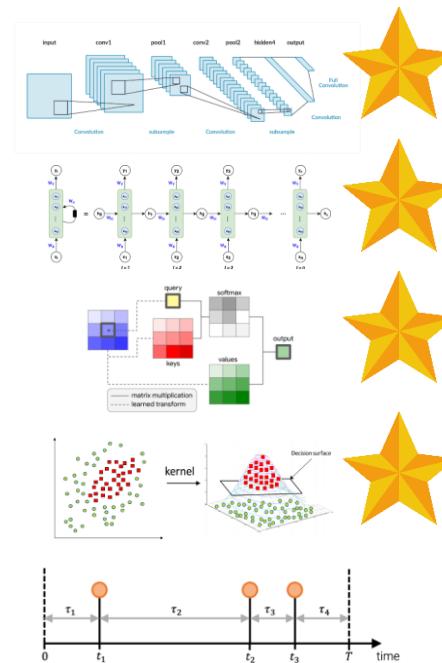
[1] Da Xu, Chuanwei Ruan, Evren Körpeoglu, Sushant Kumar, Kannan Achan: Inductive representation learning on temporal graphs. ICLR 2020

# Other Time Kernel Functions [1]

Feature maps specified by $[\phi_{2i}(t), \phi_{2i+1}(t)]$	Origin	Parameters	Interpretations of $\omega$
$[\cos(\omega_i(\mu)t), \sin(\omega_i(\mu)t)]$	Bochner's	$\mu$ : location-scale parameters specified for the <i>reparametrization trick</i> .	$\omega_i(\mu)$ : converts the $i^{th}$ sample (drawn from auxiliary distribution) to target distribution under location-scale parameter $\mu$ .
$[\cos(g_\theta(\omega_i)t), \sin(g_\theta(\omega_i)t)]$	Bochner's	$\theta$ : parameters for the inverse CDF $F^{-1} = g_\theta$ .	$\omega_i$ : the $i^{th}$ sample drawn from the auxiliary distribution.
$[\cos(\tilde{\omega}_i t), \sin(\tilde{\omega}_i t)]$	Bochner's	$\{\tilde{\omega}\}_{i=1}^d$ : transformed samples under non-parametric inverse CDF transformation.	$\tilde{\omega}_i$ : the $i^{th}$ sample of the underlying distribution $p(\omega)$ in Bochner's Theorem.
$[\sqrt{c_{2i,k}} \cos(\omega_j t), \sqrt{c_{2i+1,k}} \sin(\omega_j t)]$	Mercer's	$\{c_{i,k}\}_{i=1}^{2d}$ : the Fourier coefficients of corresponding $\mathcal{K}_{\omega_j}$ , for $j = 1, \dots, k$ .	$\omega_j$ : the frequency for kernel function $\mathcal{K}_{\omega_j}$ (can be parameters).

# In this tutorial

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process



# DyRep: Learning Representations over Dynamic Graphs [1]

- **Task:** node-level representation learning
- **Natural dynamic:** evolving graph structures and node features w.r.t time
- **Goal:** link prediction

# DyRep: Learning Representations over Dynamic Graphs [1]

- Temporal Point Process (TPP) [2]
  - A user is tweeting, they tweeted at time  $t_1 = 8:00$  am,  $t_2 = 10:00$  am,  $t_3 = 11:00$  am, what is  $t_4 = ?$



- TPP is a model that could fit the process of  $t_1$ ,  $t_2$ , and  $t_3$  to predict  $t_4$

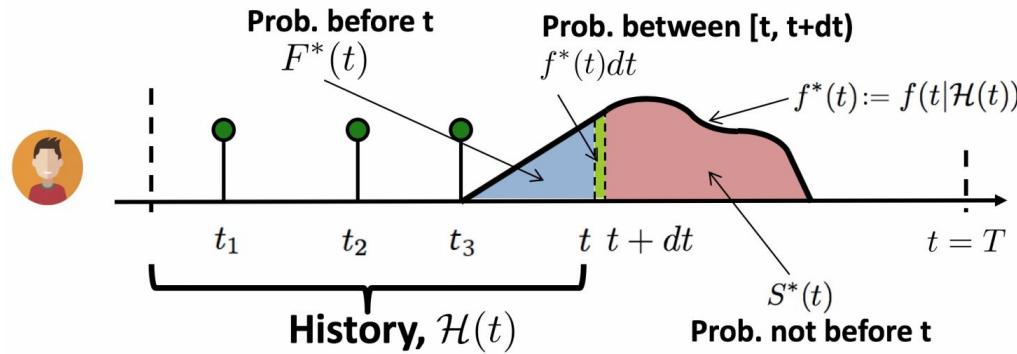
# DyRep: Learning Representations over Dynamic Graphs [1]

- Temporal Point Process (TPP) [2]
  - Given the history of events  $\mathcal{H}(t) = \{t_1, \dots, t_{i-1}\}$ , we need to model,
    - A conditional probability density function  $f^* = f(t|\mathcal{H}(t))$ , which is the conditional probability that the next event  $t$  will occur during the interval  $[t, t+dt)$
    - A cumulative distribution function  $F^*(t) = F(t|\mathcal{H}(t)) = \int_{t_{i-1}}^t f^*(\tau)d\tau$ , which is the conditional probability that the next event will occur before  $t$
    - A complementary of  $F^*(t)$ ,  $S^*(t) = S(t|\mathcal{H}(t)) = 1 - F^*(t)$ , the conditional probability that the next event will not occur before time  $t$



# DyRep: Learning Representations over Dynamic Graphs [1]

- Temporal Point Process (TPP) [2]
  - Given the history  $\mathcal{H}(t) = \{t_1, \dots, t_{i-1}\}$ , we need to model,
    - $f^*(t) = f(t|\mathcal{H}(t))$ : next event  $t$  will occur during the interval  $[t, t+dt)$
    - $F^*(t) = F(t|\mathcal{H}(t)) = \int_{t_{i-1}}^t f^*(\tau)d\tau$ : next event will occur before  $t$
    - $S^*(t) = S(t|\mathcal{H}(t)) = 1 - F^*(t)$ : next event will not occur before time  $t$

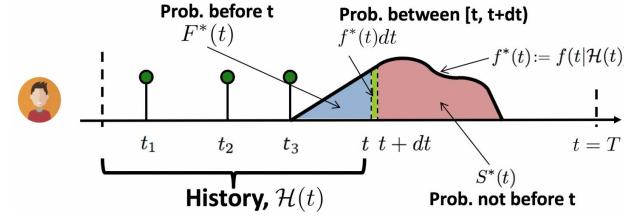


[1] Rakshit Trivedi, Mehrdad Farajtabar, Prasenjeet Biswal, Hongyuan Zha: DyRep: Learning Representations over Dynamic Graphs. ICLR 2019  
 [2] Upadhyay et al., Temporal Point Processes: <https://courses.mpi-sws.org/hcml-ws18/lectures/TPP.pdf>

# DyRep: Learning Representations over Dynamic Graphs [1]

- Temporal Point Process (TPP) [2]
  - The conditional intensity function  $\lambda^*(t) = \lambda(t|\mathcal{H}(t))$ , i.e., the conditional probability that the next event will happen during  $[t, t+dt)$ , is defined as follows

$$\lambda^*(t)dt = \frac{f^*(t)dt}{S^*(t)}$$



- $\lambda^*(t)$  can be also understood as the instantaneous rate of events per time of unit, e.g.,  $\lambda^*(t) = 10$  tweets/minute
- Using the form of  $\lambda^*(t)$  also contributes to TPP model parameterization and model reusability [2]

# DyRep: Learning Representations over Dynamic Graphs [1]

- Temporal Point Process (TPP) [2]
  - Different forms of functions model the intensity function  $\lambda^*(t)$ , e.g.,
    - Homogeneous Poisson process
      - $\lambda^*(t) = \mu \geq 0$  10 tweets/minute
    - Inhomogeneous Poisson process
      - $\lambda^*(t) = g_\theta(t) \geq 0$  2 tweets/@8:35am, 25 tweets/@2:58pm, ...
    - Hawkes process
      - $\lambda^*(t) = \mu + \alpha \sum_{t_i \in \mathcal{H}(t)} \kappa_\omega(t - t_i), \kappa_\omega(t) = \exp(-\omega t)$
  - The parameters are obtained by fitting the model with the observation and maximizing the log-likelihood

# DyRep: Learning Representations over Dynamic Graphs [1]

- Model Temporal Point Process (TPP) for Graphs
  - Each **edge connection** is considering as an **event**  $(u, v, t)$
  - We want to predict whether a node  $u$  and a node  $v$  will connect at time  $t$ , given node  $u$ 's history and node  $v$ 's history before  $t$

# DyRep: Learning Representations over Dynamic Graphs [1]

- Model Temporal Point Process (TPP) for Graphs
  - Intensity functions for graphs, i.e., an edge connection between nodes  $u$  and  $v$

$$\lambda^{u,v}(t) = f(g^{u,v}(\bar{t}))$$

$\bar{t}$  means the timestamp just before the current event

# DyRep: Learning Representations over Dynamic Graphs [1]

- Model Temporal Point Process (TPP) for Graphs
  - Intensity functions for graphs, i.e., an edge connection between nodes  $u$  and  $v$

the outer function  $f()$  is a softplus function  
with trainable a parameter to  
make sure positive output

$$\lambda^{u,v}(t) = f(g^{u,v}(\bar{t}))$$

$$f(x) = \psi \log(1 + \exp(x/\psi))$$

$$g^{u,v}(\bar{t}) = \omega [z^u(\bar{t}); z^v(\bar{t})]$$

the inner function  $g()$  computes  
the compatibility of the catenation

- Now, the question is how to get the node embeddings, e.g.,  $z^v(\bar{t})$ ?

# DyRep: Learning Representations over Dynamic Graphs [1]

- Model Temporal Point Process (TPP) for Graphs
  - How to get node embeddings, e.g.,  $\mathbf{z}^v(\bar{t})$ ?
    - **Self-Propagation**: w.r.t its historical behavior
    - **Exogeneous Drive**: for the smooth update of the current
    - **Localized Embedding Propagation**: message passing within second-order proximity
  - Suppose node  $u$  and node  $v$  participating in any type of event at time  $t$ 
    - E.g., for the  $p$ -th event of node  $v$  at time  $t$

$$\mathbf{z}^v(t_p) = \sigma\left( \underbrace{\mathbf{W}^{struct} \mathbf{h}_{struct}^u(\bar{t}_p)}_{\text{Localized Embedding Propagation}} + \underbrace{\mathbf{W}^{rec} \mathbf{z}^v(\bar{t}_p^v)}_{\text{Self-Propagation}} + \underbrace{\mathbf{W}^t(t_p - \bar{t}_p^v)}_{\text{Exogenous Drive}} \right)$$

# DyRep: Learning Representations over Dynamic Graphs [1]

- Suppose node  $u$  and node  $v$  participating in any type of event at time  $t$ 
  - For the  $p$ -th event of node  $v$  at time  $t$

$$\mathbf{z}^v(t_p) = \sigma(\underbrace{\mathbf{W}^{struct} \mathbf{h}_{struct}^u(\bar{t}_p)}_{\text{Localized Embedding Propagation}} + \underbrace{\mathbf{W}^{rec} \mathbf{z}^v(\bar{t}_p)}_{\text{Self-Propagation}} + \underbrace{\mathbf{W}^t(t_p - \bar{t}_p)}_{\text{Exogenous Drive}})$$

$\mathbf{h}_{struc}^u$  is the aggregation on node  $u$ 's neighbors

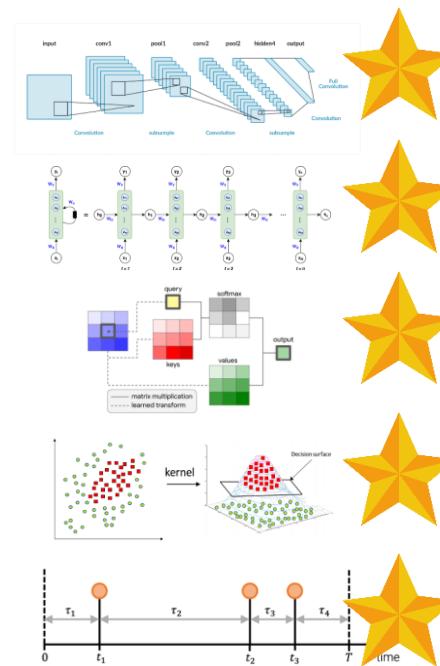
$$h_{struc}^u(\bar{t}) = \max(\{\sigma(q_{ui}(t) \cdot h^i(\bar{t})) \mid \forall i \in N_u(\bar{t})\})$$

$q_{ui}(t)$  can be understood as the weight of the connection

$$h^i(\bar{t}) = \mathbf{W}z^i(\bar{t}) + \mathbf{b}$$

# In this tutorial

- Covered works for natural dynamics in GNNs
  - Temporal GNNs with Convolutional Operations
  - Temporal GNNs with Recurrent Units
  - Temporal GNNs with Time Attention
  - Temporal GNNs with Time Kernel
  - Temporal GNNs with Temporal Point Process



# There are also many great research works on natural dynamics

- Ma et al. **Streaming Graph Neural Networks**. SIGIR 2020
- Rossi et al. **Temporal Graph Networks for Deep Learning on Dynamic Graphs**. CoRR (2020)
- Tian et al. **Self-supervised Representation Learning on Dynamic Graphs**. CIKM 2021
- Fu et al. **SDG: A Simplified and Dynamic Graph Neural Network**. SIGIR 2021
- You et al. **ROLAND: Graph Learning Framework for Dynamic Graphs**. KDD 2022
- Cong et al. **Do We Really Need Complicated Model Architectures For Temporal Networks?** ICLR 2023
- Many more.....

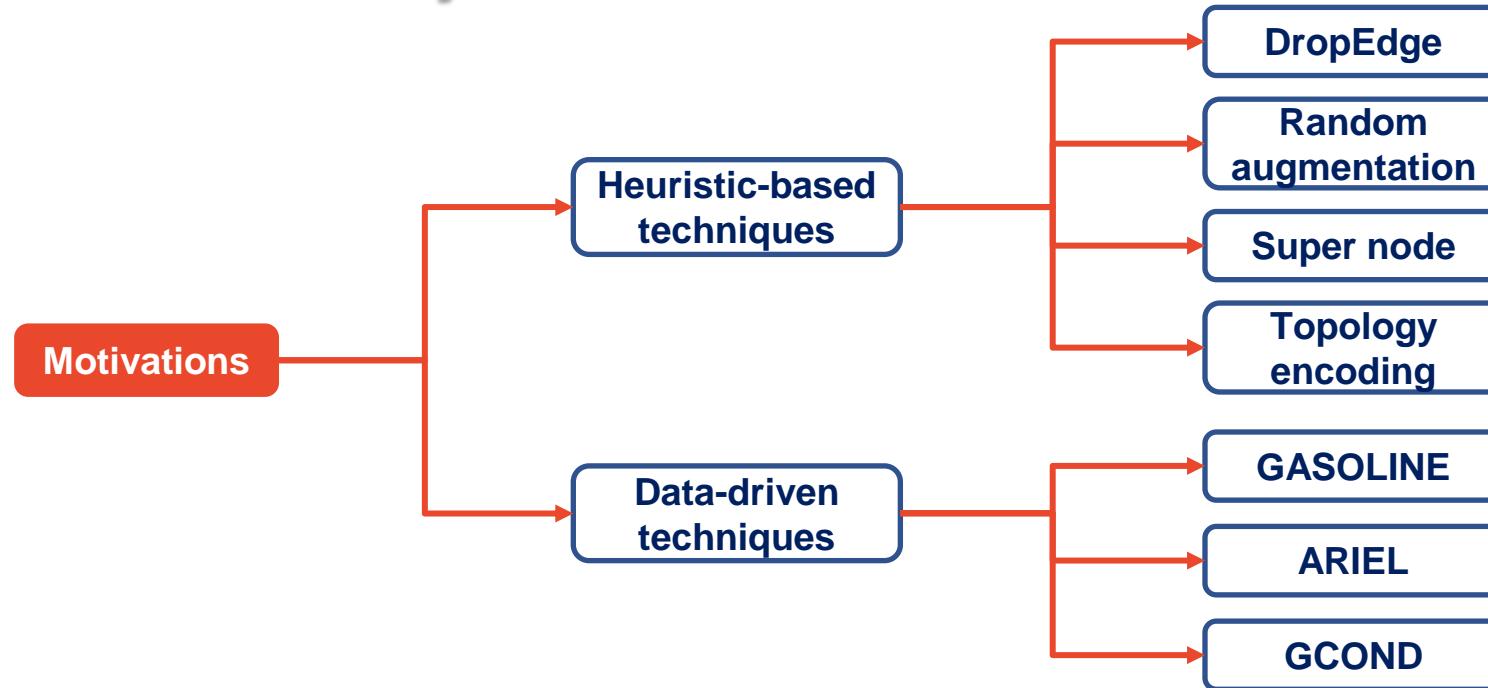
# Contents

- Part I – Introduction 
- Part II – Natural Dynamics in GNNs 
- 30 mins Coffee Break
- Part III – Artificial Dynamics in GNNs
- Part IV – Challenges and Future Directions
- Q&A

# Contents

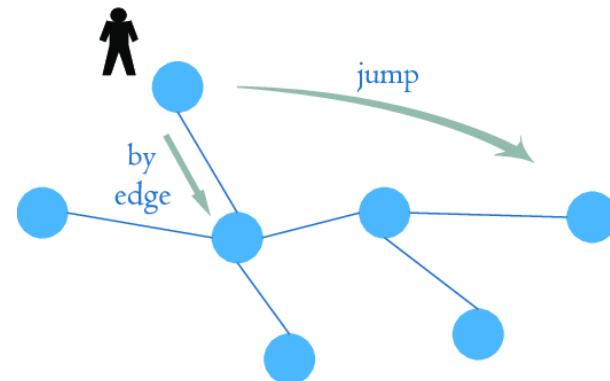
- Part I – Introduction 
- Part II – Natural Dynamics in GNNs 
- 30 mins Coffee Break 
- Part III – Artificial Dynamics in GNNs
- Part IV – Challenges and Future Directions
- Q&A

# Artificial dynamics in GNNs



# Artificial dynamics in GNNs – definition

- “**Artificial jump**” [1] is proposed to adjust the graph topology for PageRank realizing the personal ranking function on graphs.
- Artificial dynamics in the **GNN** context: modifying the given graph(s) for specific functionality.
- A.k.a. graph data augmentation (for GNNs) [2,3].



# Artificial dynamics in GNNs – motivation (1)

The given graph(s) could be **noisy and poisoned**.

Example 1: An academic literature classification system



# Artificial dynamics in GNNs – motivation (2)

The given graph(s) could be **not informative**.

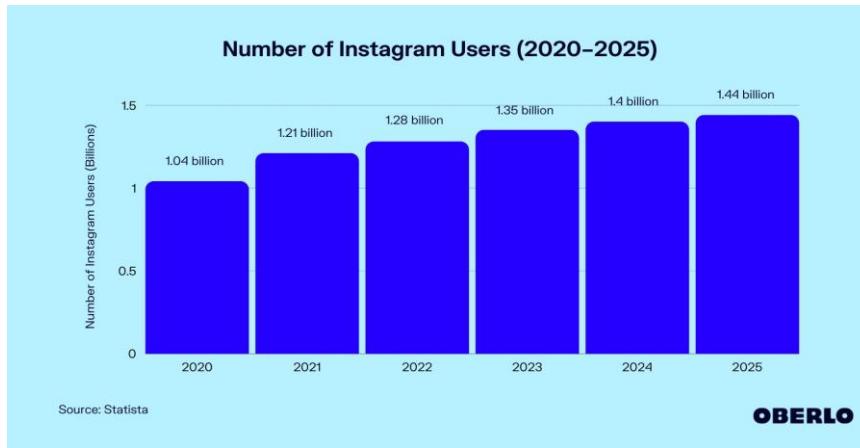
Example 2: Friend recommendation on social networks



# Artificial dynamics in GNNs – motivation (3)

The given graph(s) could be **huge**.

Example 3: representation learning on real-world networks

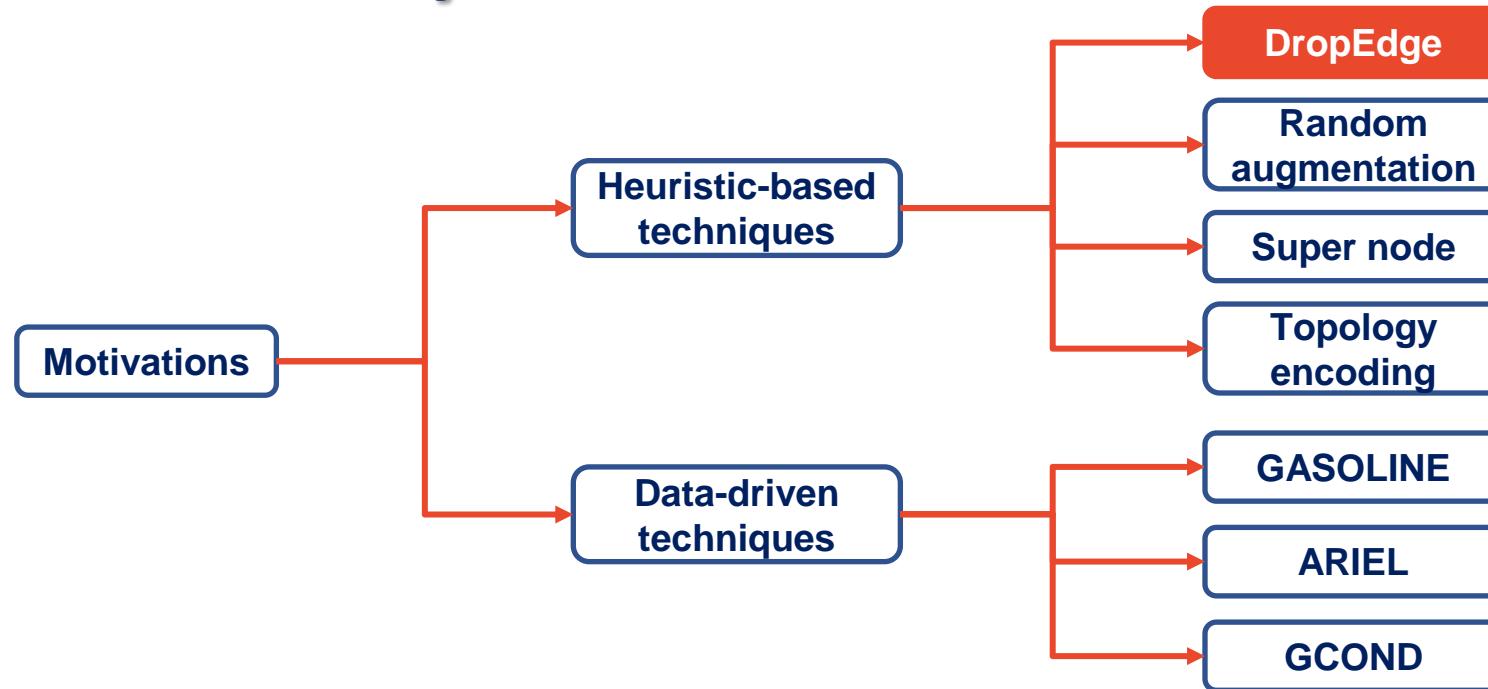


Task category	Name	Package	#Graphs	#Total nodes	#Total edges
Node-level	<a href="#">MAG240M</a>	<code>&gt;=1.3.2</code>	1	244,160,499	1,728,364,232
Link-level	<a href="#">WikiKG90Mv2</a>	<code>&gt;=1.3.3</code>	1	91,230,610	601,062,811
Graph-level	<a href="#">PCQM4Mv2</a>	<code>&gt;=1.3.2</code>	3,746,619	52,970,652	54,546,813

Figure credit: [1] [2]

And many more ...

# Artificial dynamics in GNNs



# DropEdge: Towards deep graph convolutional networks on node classification

**Task:** node classification

**Artificial dynamic:** randomly removes a certain number of edges from the input graph at each training epoch

**Goal:** reduce the message passing of GNN to alleviate over-smoothing

$$\mathbf{A}_{drop} = \mathbf{A} - \mathbf{A}'$$

# DropEdge: alleviate over-smoothing

- **PageRank's convergence**,  $\tilde{\mathbf{A}}$  is a normalized adjacency matrix,  $\pi$  is only determined by  $\tilde{\mathbf{A}}$  for all the  $\mathbf{x}$  except 0:

$$\lim_{k \rightarrow \infty} \tilde{\mathbf{A}}^k \mathbf{x} = \pi$$

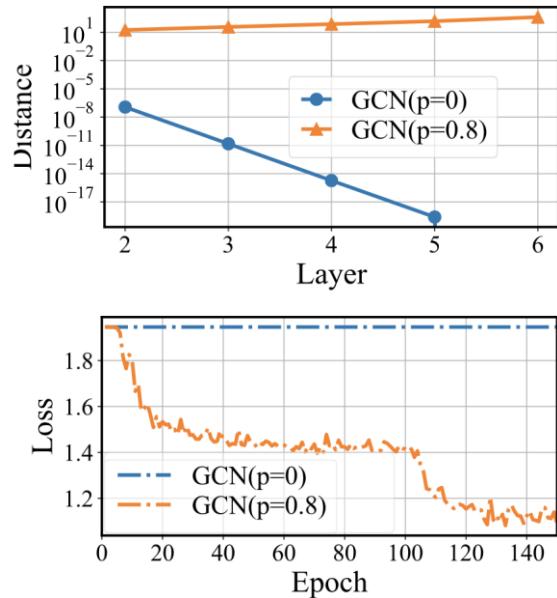
- The following two GNNs share **similar spectral characteristics** and **empirical performance**:

- **GCN**:  $\sigma_2(\tilde{\mathbf{A}}\sigma_1(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W}_1)\mathbf{W}_2)$
- **SGC**:  $\sigma_2(\tilde{\mathbf{A}}\tilde{\mathbf{A}}\mathbf{X}\mathbf{W}_1\mathbf{W}_2) = \sigma(\tilde{\mathbf{A}}^2\mathbf{X}\mathbf{W})$

- **SGC with many layers**:  $\lim_{k \rightarrow \infty} \sigma(\tilde{\mathbf{A}}^k \mathbf{X}\mathbf{W})$ , information from initial node feature  $\mathbf{X}$  is missing ( $\tilde{\mathbf{A}}^k \mathbf{X}\mathbf{W}$  will converge to a matrix **independent** of  $\mathbf{X}\mathbf{W}$ ).
- **Dropedge**:  $\mathbf{A}_{drop} = \mathbf{A} - \mathbf{A}'$

# Empirical results of DropEdge

(1) Embeddings are not over-smoothed and training moves forward:

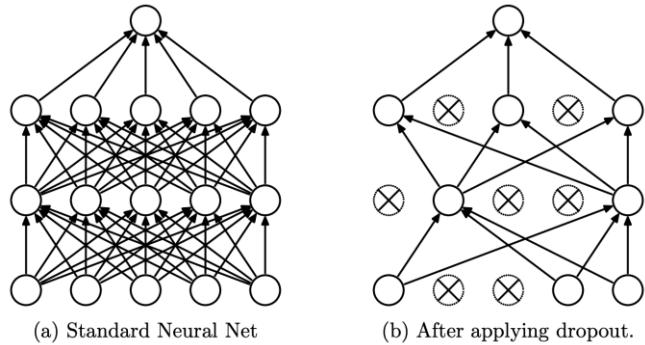


(2) Better performance, especially for the deep GNN cases

Dataset	Backbone	2 layers		8 layers		32 layers	
		Original	DropEdge	Original	DropEdge	Original	DropEdge
Cora	GCN	86.10	<b>86.50</b>	78.70	<b>85.80</b>	71.60	<b>74.60</b>
	ResGCN	-	-	85.40	<b>86.90</b>	85.10	<b>86.80</b>
	JKNet	-	-	86.70	<b>87.80</b>	87.10	<b>87.60</b>
	IncepGCN	-	-	86.70	<b>88.20</b>	87.40	<b>87.70</b>
	GraphSAGE	87.80	<b>88.10</b>	84.30	<b>87.10</b>	31.90	<b>32.20</b>
Citeseer	GCN	75.90	<b>78.70</b>	74.60	<b>77.20</b>	59.20	<b>61.40</b>
	ResGCN	-	-	77.80	<b>78.80</b>	74.40	<b>77.90</b>
	JKNet	-	-	79.20	<b>80.20</b>	71.70	<b>80.00</b>
	IncepGCN	-	-	79.60	<b>80.50</b>	72.60	<b>80.30</b>
	GraphSAGE	78.40	<b>80.00</b>	74.10	<b>77.10</b>	37.00	<b>53.60</b>
Pubmed	GCN	90.20	<b>91.20</b>	90.10	<b>90.90</b>	84.60	<b>86.20</b>
	ResGCN	-	-	89.60	<b>90.50</b>	90.20	<b>91.10</b>
	JKNet	-	-	90.60	<b>91.20</b>	89.20	<b>91.30</b>
	IncepGCN	-	-	90.20	<b>91.50</b>	OOM	<b>90.50</b>
	GraphSAGE	90.10	<b>90.70</b>	90.20	<b>91.70</b>	41.30	<b>47.90</b>
Reddit	GCN	96.11	<b>96.13</b>	96.17	<b>96.48</b>	45.55	<b>50.51</b>
	ResGCN	-	-	96.37	<b>96.46</b>	93.93	<b>94.27</b>
	JKNet	-	-	96.82	<b>97.02</b>	OOM	OOM
	IncepGCN	-	-	96.43	<b>96.87</b>	OOM	OOM
	GraphSAGE	96.22	<b>96.28</b>	96.38	<b>96.42</b>	96.43	<b>96.47</b>

# Related works of DropEdge

- **Dropout [1]:** randomly zeroing a part of the feature transformation matrices

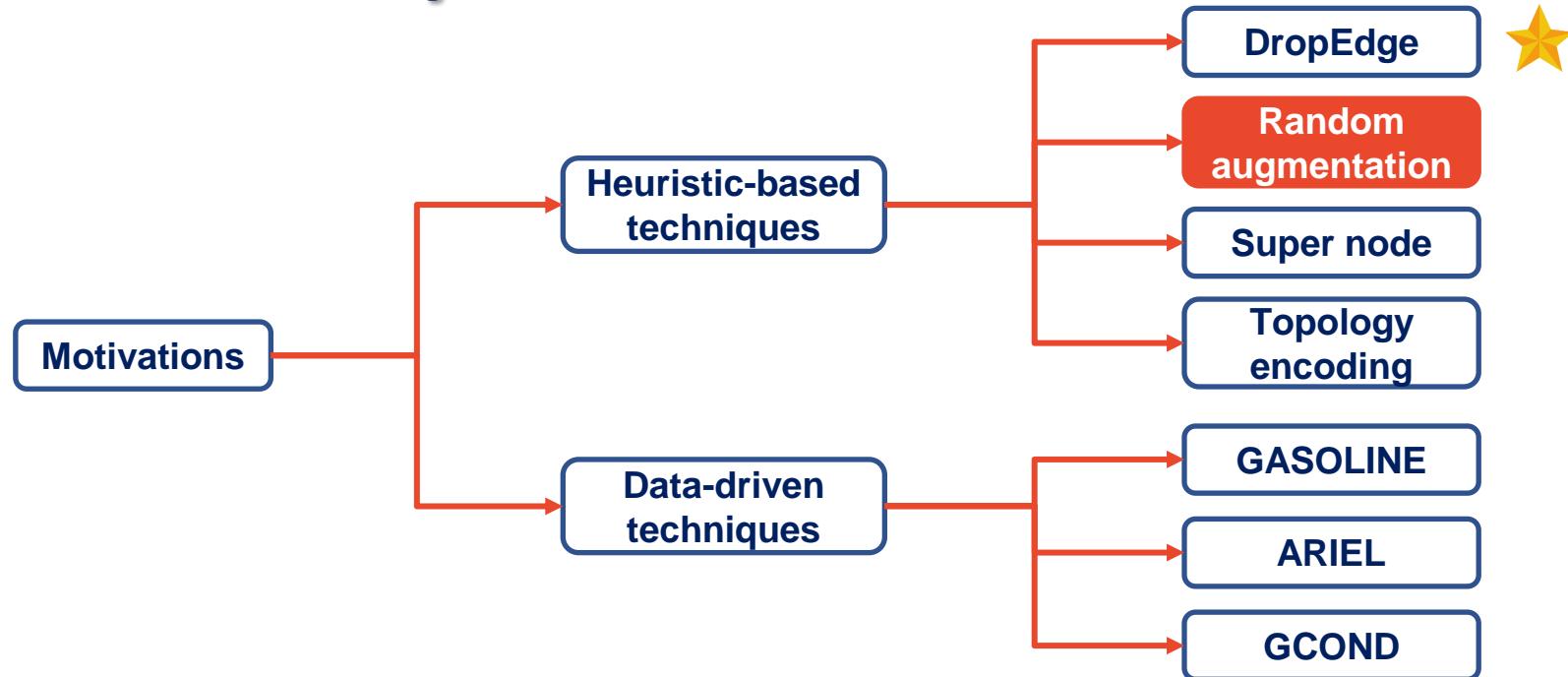


$$\sigma(\tilde{\mathbf{A}}\mathbf{X}\mathbf{W})$$

Target of DropEdge      Target of Dropout

- **Sampling-based methods:** GraphSAGE [2], FastGCN [3], GraphSAINT[4], and many more
- **Other applications:** (1) improving GNN scalability, (2) alleviating over-fitting, (3) generating different views of the given graphs, and many more

# Artificial dynamics in GNNs

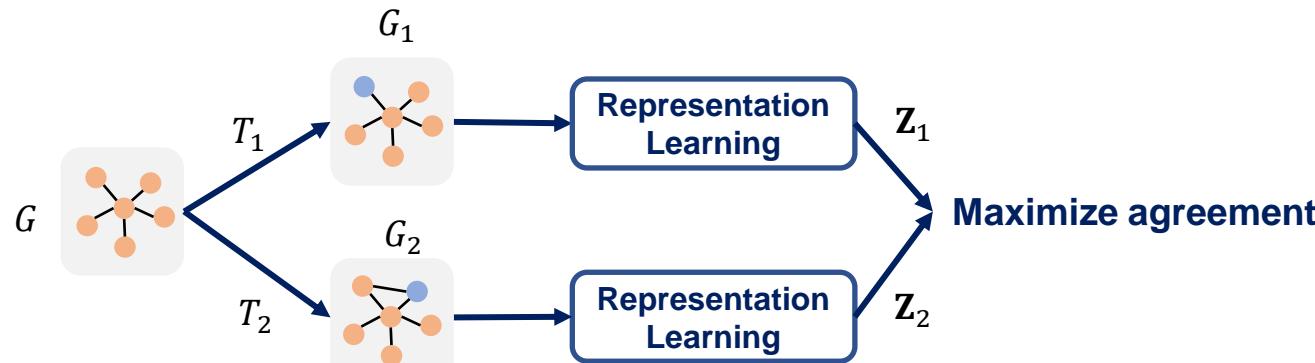


# Graph contrastive learning (GraphCL) with augmentations

**Task:** unsupervised **graph-level** representation learning

**Artificial dynamic:** using random augmentations to generate views

**Goal:** contrastive learning on graph data



$T_1, T_2$  : Heuristic-based augmentations

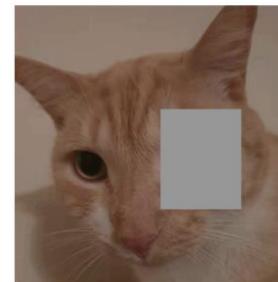
# Heuristic-based augmentations from GraphCL

Image data [1]:

Original image



Random erasing



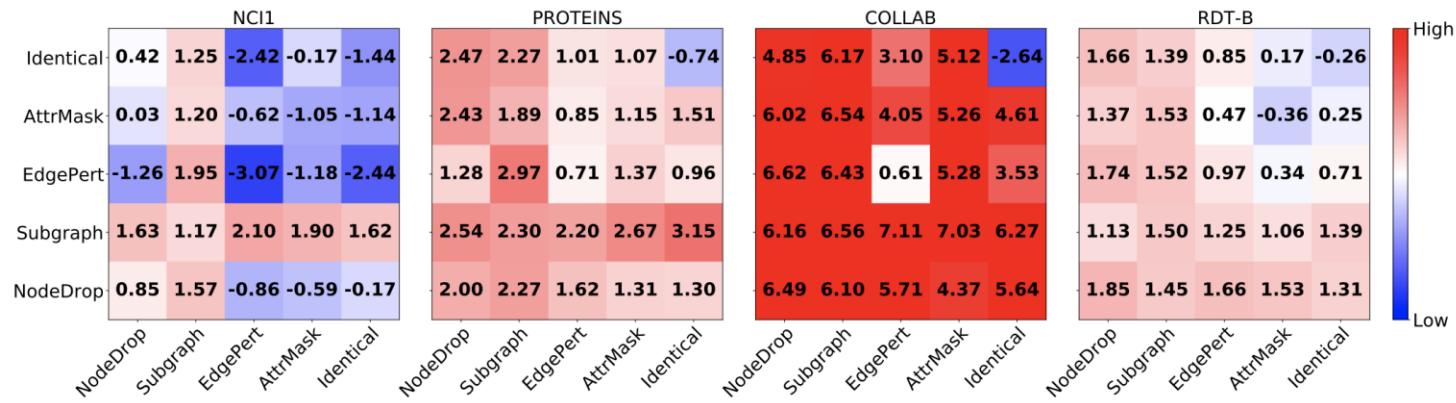
Graph data [2]:

**Table 1:** Overview of data augmentations for graphs.

Data augmentation	Type	Underlying Prior
Node dropping	Nodes, edges	Vertex missing does not alter semantics.
Edge perturbation	Edges	Semantic robustness against connectivity variations.
Attribute masking	Nodes	Semantic robustness against losing partial attributes per node.
Subgraph	Nodes, edges	Local structure can hint the full semantics.

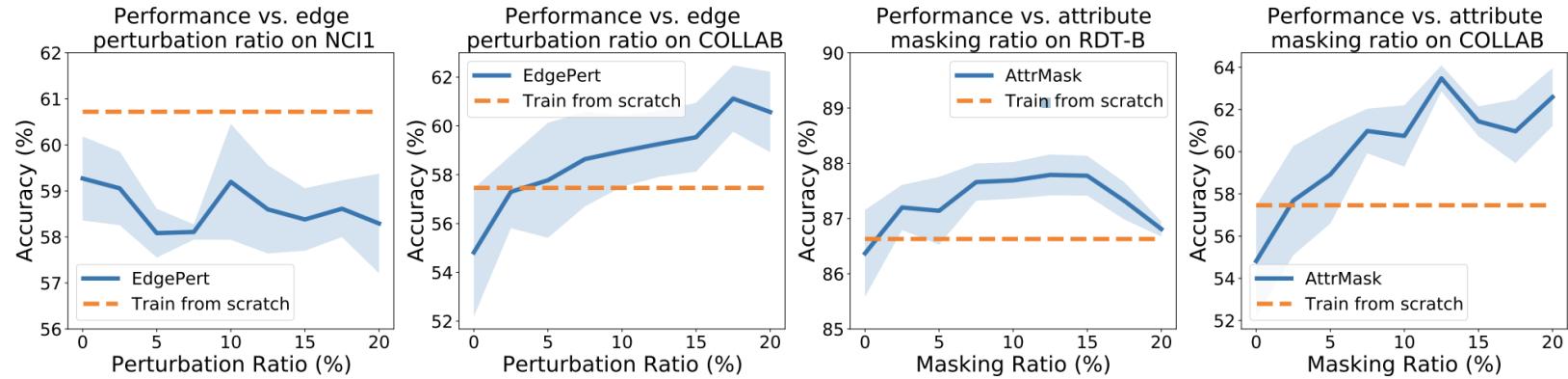
# Empirical results of GraphCL in terms of augmentations

- The **selection** of augmentation methods is critical
- The augmentations **cannot generalize** across different datasets



**Figure 2:** Semi-supervised learning accuracy gain (%) when contrasting different augmentation pairs, compared to training from scratch, under four datasets: NCI1, PROTEINS, COLLAB, and RDT-B. Pairing “Identical” stands for a no-augmentation baseline for contrastive learning, where the positive pair diminishes and the negative pair consists of two non-augmented graphs. Warmer colors indicate better performance gains. The baseline training-from-scratch accuracies are 60.72%, 70.40%, 57.46%, 86.63% for the four datasets respectively.

# Trials and errors of heuristic-based augmentations



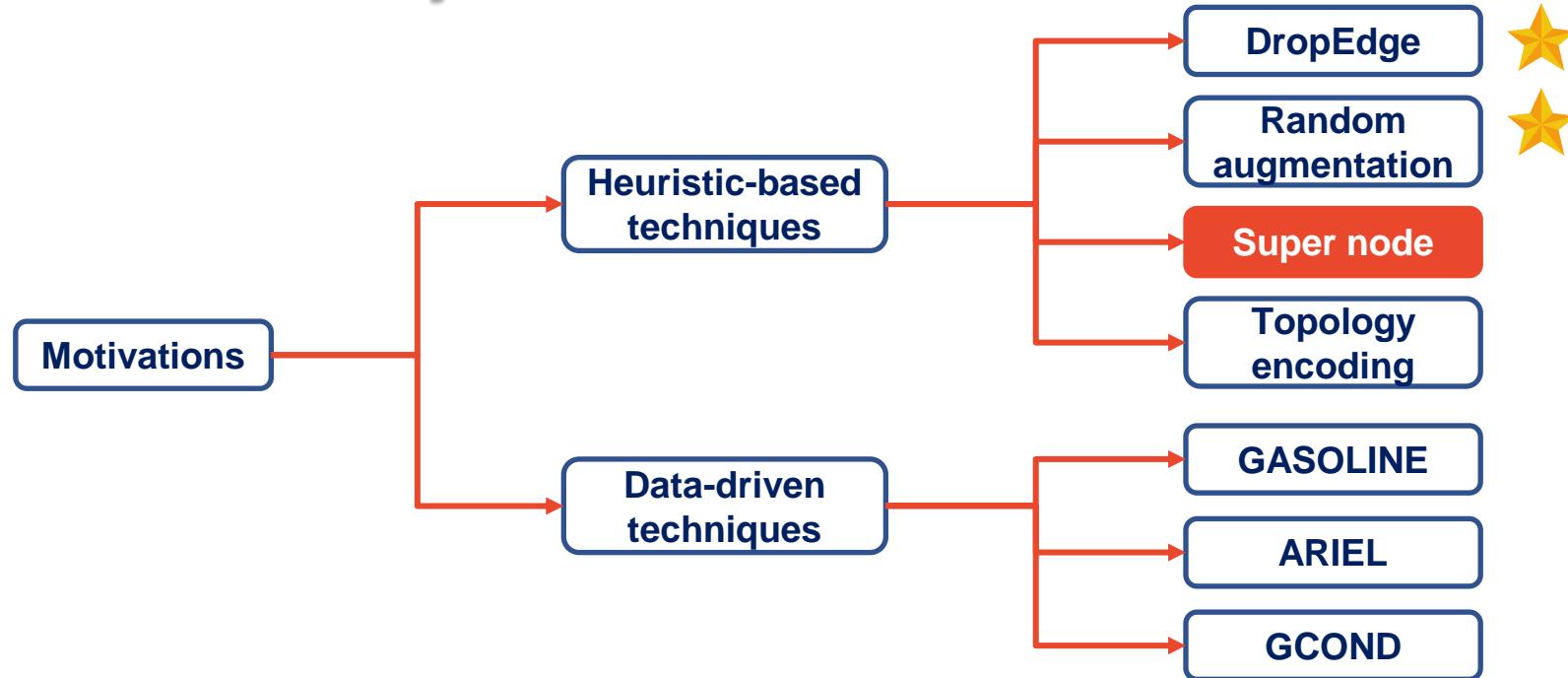
**Figure 4:** Performance versus augmentation strength. Left two figures implemented edge perturbation with different ratios. The right two figures apply attribute masking with different masking ratios.

The selection of augmentation tools and hyperparameters is **critical** but **labor-intensive**.

# Related works of GraphCL

- Most of the augmentations adopted by GraphCL [1] are **heuristic**, e.g., uniformly random dropping nodes.
- **Learn to select** the heuristic-based augmentation combinations [2].
- **Adapt the probability of heuristic-based augmentation** according to some metrics [3]. E.g., adapt the node dropping probability according to the node centrality.
- Many more ...

# Artificial dynamics in GNNs



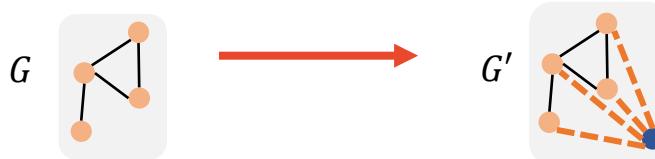
# Inserting super node

**Representation works:** [1-4]

**Task:** node/graph property prediction

**Artificial dynamic:** inserting a super node (a.k.a., virtual/master/dummy node) into the given graph

**Goal:** (1) to improve the message passing of GNN, and (2) can be used as a readout node for graph-level tasks.



[1] Gilmer, Justin, et al. "Neural message passing for quantum chemistry." International conference on machine learning. PMLR, 2017.

[2] J. Li, D. Cai, and X. He. Learning graph-level representation for drug discovery. arXiv preprint arXiv:1709.03741, 2017

[3] K. Ishiguro, S.-i. Maeda, and M. Koyama. Graph warp module: an auxiliary module for boosting the power of graph neural networks in molecular graph analysis. arXiv preprint arXiv:1902.01020, 2019.

[4] T. Pham, T. Tran, H. Dam, and S. Venkatesh. Graph classification via deep learning with virtual nodes. arXiv preprint arXiv:1708.04357, 2017

# Details comparison of inserting super node

	[1]	[2]	[3]	[4]
<b>Weighted?</b>	No	No	Yes. Learnable gates.	No
<b>Directed?</b>	No	Yes	No	No
<b>Virtual node feature</b>	Extra node features	Zeros at initialization	Graph features. E.g., graph diameter	Graph features
<b>As readout?</b>	No	Yes	No	Yes



I [1] Gilmer, Justin, et al. "Neural message passing for quantum chemistry." International conference on machine learning. PMLR, 2017.

[2] J. Li, D. Cai, and X. He. Learning graph-level representation for drug discovery. arXiv preprint arXiv:1709.03741, 2017

[3] K. Ishiguro, S.-i. Maeda, and M. Koyama. Graph warp module: an auxiliary module for boosting the power of graph neural networks in molecular graph analysis. arXiv preprint arXiv:1902.01020, 2019.

[4] T. Pham, T. Tran, H. Dam, and S. Venkatesh. Graph classification via deep learning with virtual nodes. arXiv preprint arXiv:1708.04357, 2017

# Effectiveness of inserting super node on OGB dataset

Statistics from the leaderboard of the OGB datasets [1,2].

On ogbg-molhiv dataset

Rank	Method	Test ROC- AUC	Validation ROC-AUC
25	GIN+ <b>virtual node</b> +FLAG	0.7748 ± 0.0096	0.8438 ± 0.0128
27	GIN+ <b>virtual node</b>	0.7707 ± 0.0149	0.8479 ± 0.0068
31	GCN+ <b>virtual node</b>	0.7599 ± 0.0119	0.8384 ± 0.0091

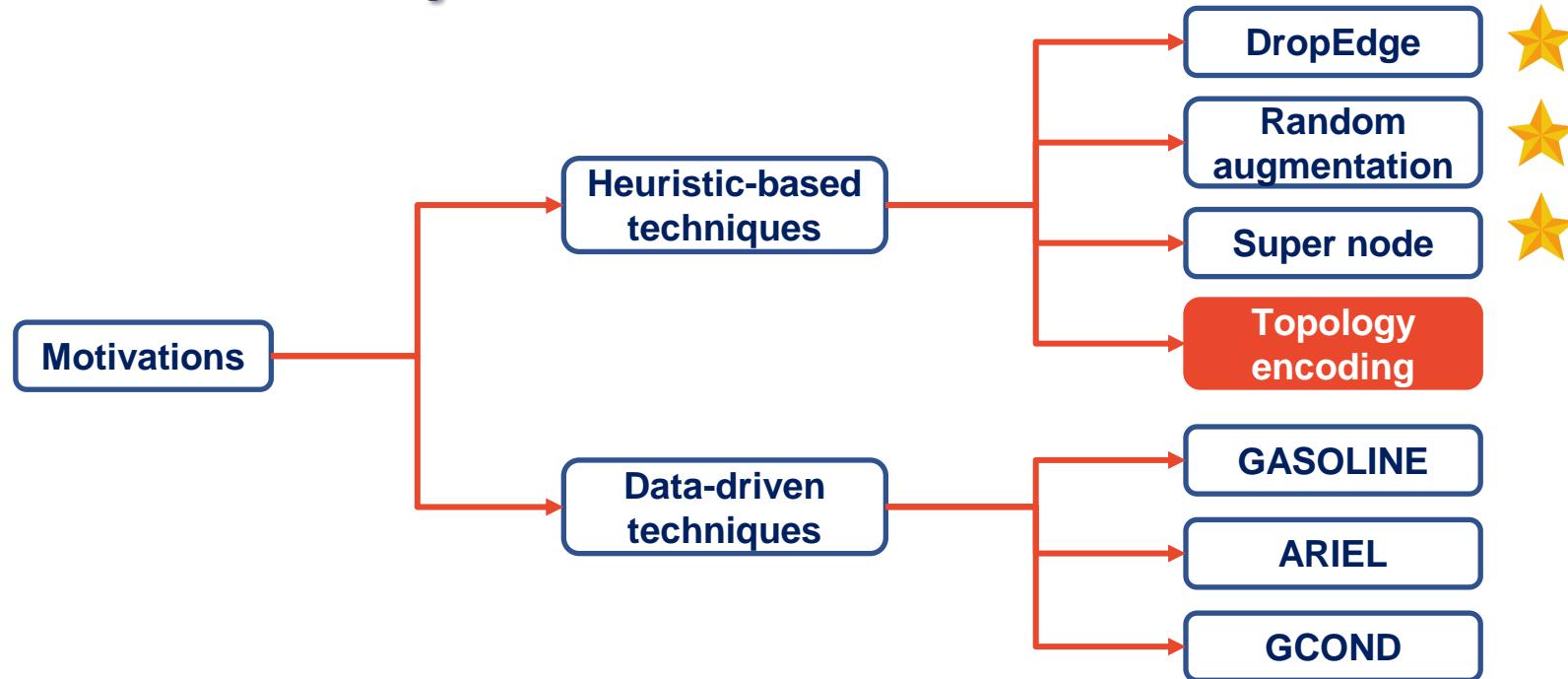
On ogbg-molhiv dataset

Rank	Method	Test AP	Validation AP
4	NestedGIN+ <b>vi rtual node</b>	0.3007 ± 0.0037	0.3059 ± 0.0056
10	GINE+ <b>virtual nodes</b>	0.2917 ± 0.0015	0.3065 ± 0.0030
13	DeeperGCN+ <b>virtual node</b> +FLAG	0.2842 ± 0.0043	0.2952 ± 0.0029

[1] Hu, Weihua, et al. "Open graph benchmark: Datasets for machine learning on graphs." Advances in neural information processing systems 33 (2020): 22118-22133.

[2] [https://ogb.stanford.edu/docs/leader\\_graphprop/](https://ogb.stanford.edu/docs/leader_graphprop/)

# Artificial dynamics in GNNs



# Distance encoding: Design provably more powerful neural networks for graph representation learning

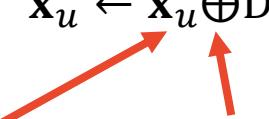
**Task:** representation learning of a target (sub)graph  $S$

**Artificial dynamic:** generating extra node features by encoding the topology information (distance in this paper) which is hard to be caught by 1-WL GNNs\*.

**Goal:** to enrich the informativeness of node features

$$\mathbf{x}_u \leftarrow \mathbf{x}_u \oplus \text{DE}(u, S)$$

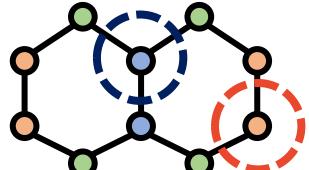
Raw feature of node  $u$     concatenate



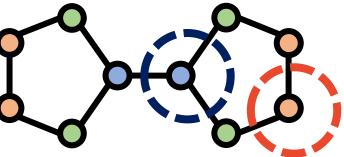
1-WL GNNs refer to GNNs whose expressiveness is limited by the 1-WL test.

# Expressiveness limitation of 1-WL GNNs

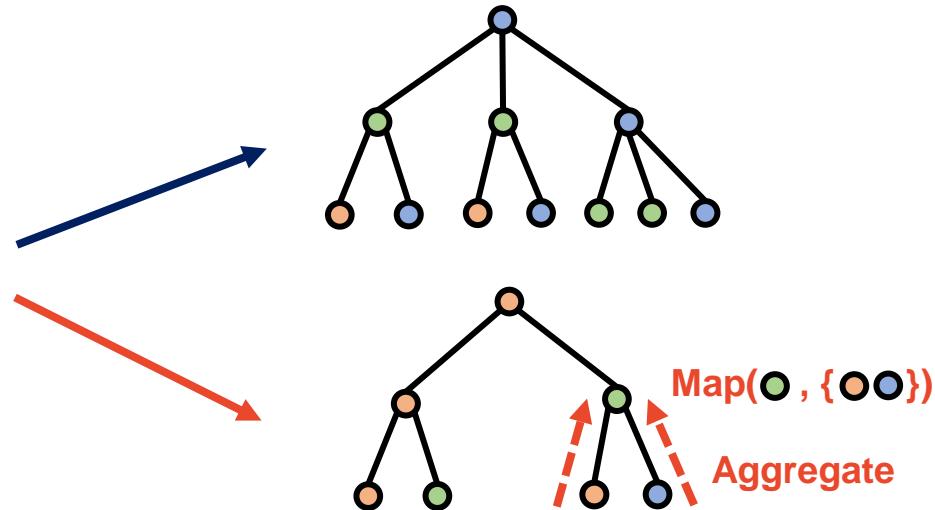
Figure credit: [1,2]



Graph 1



Graph 2



GNNs are not substructure-aware [3,4] and position-aware [5]

Node pairs from two graphs share the **same rooted subtrees**, which are used by Weisfeiler-Lehman (WL) graph isomorphism test and (most) GNN operations.

[1] Ding, Kaize, et al. "Data augmentation for deep graph learning: A survey." ACM SIGKDD Explorations Newsletter 24.2 (2022): 61-77.

[2] Sato, Ryoma. "A survey on the expressive power of graph neural networks." arXiv preprint arXiv:2003.04078 (2020).

[3] Chen, Zhengdao, et al. "Can graph neural networks count substructures?." Advances in neural information processing systems 33 (2020): 10383-10395.

[4] J. You, J. M. Gomes-Selman, R. Ying, and J. Leskovec. Identity-aware graph neural networks. In AAAI, 2021

[5] You, Jiaxuan, Rex Ying, and Jure Leskovec. "Position-aware graph neural networks." International conference on machine learning. PMLR, 2019.

# Details of the distance encoding

$$\mathbf{x}_u = \mathbf{x}_u \oplus \text{DE}(u, S)$$

$$\text{DE}(u, S) = \text{AGG}(\{\text{DE}(u, v) | v \in S\})$$



Set aggregation, e.g., SUM

DE Option 1: (generalized) PageRank score

$$\text{DE}(u, v) = \sum_{k \geq 1} \gamma_k \tilde{\mathbf{A}}^k[u, v]$$

DE Option 2: shortest path distance

$$\text{DE}(u, v) = \text{spd}(u, v)$$

# Empirical results of the distance encoding

- The size of the target subgraph  $S$  determines whether the task is node-level, link-level, or triangle-level.
- Their experiments focus on structure representation, so all the datasets are non-attributed.

Method	Data	Nodes (Task 1): Average Accuracy			Node-pairs (Task 2): AUC			Node-triads (Task 3): AUC	
		Bra.-Airports	Eur.-Airports	USA-Airports	C.elegans	NS	PB	C.elegans	NS
GCN [20]		64.55±4.18	54.83±2.69	56.58±1.11	74.03±0.99	74.21±1.72	89.78±0.99	80.94±0.51	81.72±1.50
SAGE [21]		70.65±5.33	56.29±3.21	50.85±2.83	73.91±0.32	79.96±1.40	90.23±0.74	84.72±0.40	84.06±1.14
GIN [16]		71.89±3.60 <sup>†</sup>	57.05±4.08	58.87±2.12	75.58±0.59	87.75±0.56	91.11±0.52	86.42±1.12 <sup>†</sup>	94.59±0.66 <sup>†</sup>
Struc2vec [5]		70.88±4.26	57.94±4.01 <sup>†</sup>	61.92±2.61 <sup>†</sup>	72.11±0.31	82.76±0.59	90.47±0.60	77.72±0.58	81.93±0.61
PGNN [10]		N/A	N/A	N/A	78.20±0.33	94.88±0.77	89.72±0.32	86.36±0.74	79.36±1.49
SEAL [9]		N/A	N/A	N/A	88.26±0.56 <sup>†</sup>	98.55±0.32 <sup>†</sup>	94.18±0.57 <sup>†</sup>	N/A	N/A
DE-GNN-SPD		<b>73.28±2.47</b>	56.98±2.79	<b>63.10±0.68*</b>	<b>89.37±0.17*</b>	<b>99.09±0.79</b>	<b>94.95±0.37*</b>	<b>92.17±0.72*</b>	<b>99.65±0.40*</b>
DE-GNN-LP		<b>75.10±3.80*</b>	58.41±3.20*	<b>64.16±1.70*</b>	86.27±0.33	98.01±0.55	91.45±0.41	86.24±0.18	<b>99.31±0.12*</b>
DEA-GNN-SPD		<b>75.37±3.25*</b>	57.99±2.39*	<b>63.28±1.59</b>	<b>90.05±0.26*</b>	<b>99.43±0.63*</b>	<b>94.49±0.24*</b>	<b>93.35±0.65*</b>	<b>99.84±0.14*</b>

Table 1: Performance in Average Accuracy and Area Under the ROC Curve (AUC) (mean in percentage  $\pm$  95% confidence level). <sup>†</sup> highlights the best baselines. \*, **bold font**, **bold font**\* respectively highlights the case where our models' performance exceeds the best baseline on average, by 70% confidence, by 95% confidence.

# Related works of distance encoding

- Distance [1] is not the only information that is hard to be caught by 1-WL GNNs
- GSN [2] and ID-GNN [3] augment the node feature with **the counts of various motifs (e.g., cycles)**.
- Adding **random features** [4] into the existing node features can help GNN discriminate triangle structure

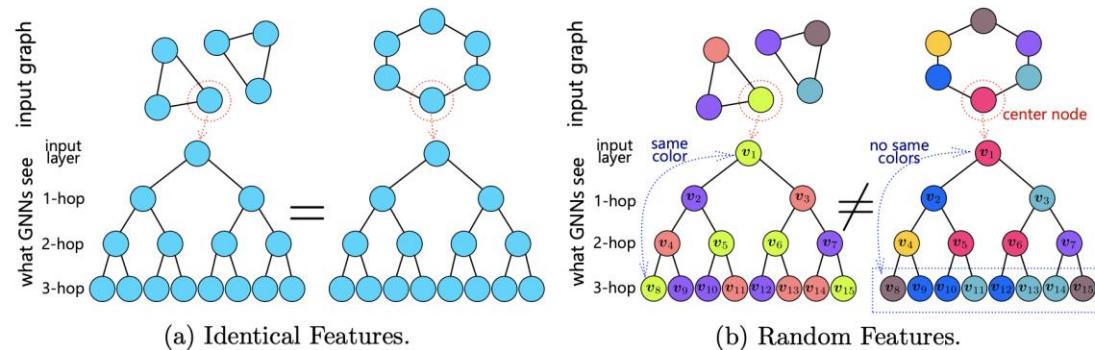
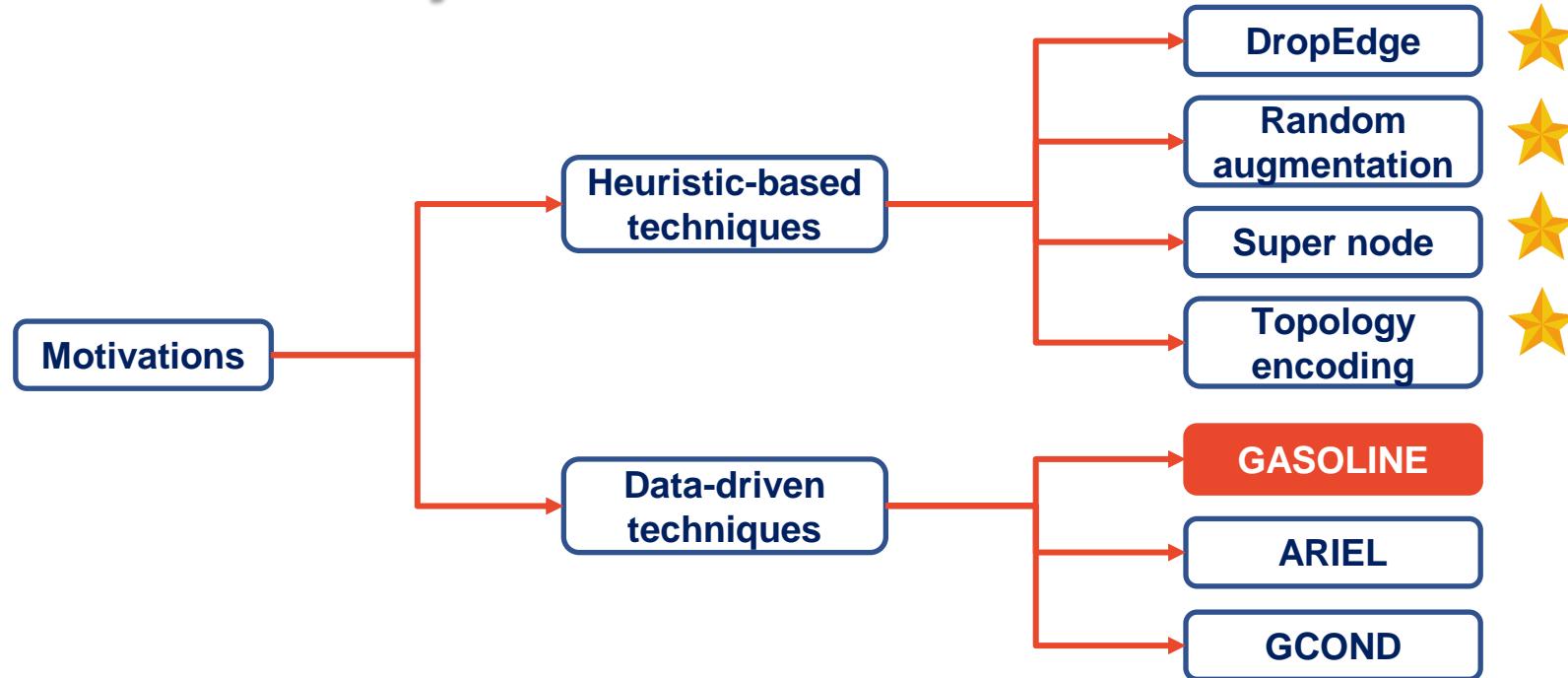


Figure credit: [4]

# Other heuristic-based artificial dynamics

- Graph diffusions [1]
- Low-rank approximated adjacency matrix [2]
- Recover missing features based on homophily assumption [3]
- Many more [4] ...

# Artificial dynamics in GNNs

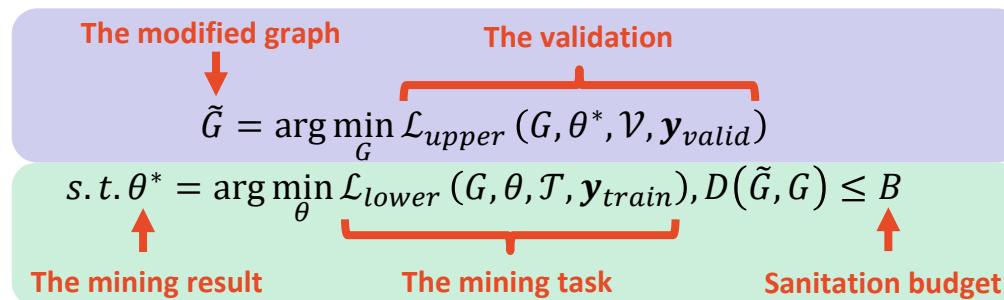


# Graph sanitation with application to node classification (GASOLINE)

**Task:** Node classification

**Artificial dynamic:** View the graph data itself as a hyper-parameter and optimize it by a meta-learning fashioned bilevel optimization objective

**Goal:** improve the node classification performance



# Instantiation of GASOLINE on node classification tasks

$$\begin{aligned}\tilde{G} &= \arg \min_G \mathcal{L}_{upper}(G, \theta^*, \mathcal{V}, \mathbf{y}_{valid}) \\ s.t. \theta^* \\ &= \arg \min_{\theta} \mathcal{L}_{lower}(G, \theta, \mathcal{T}, \mathbf{y}_{train})\end{aligned}$$

$$\tilde{G} = \arg \min_G - \sum_{i \in \mathcal{V}} \sum_{j=1}^c y_{ij} \ln f(G, \theta^*)[i, j]$$

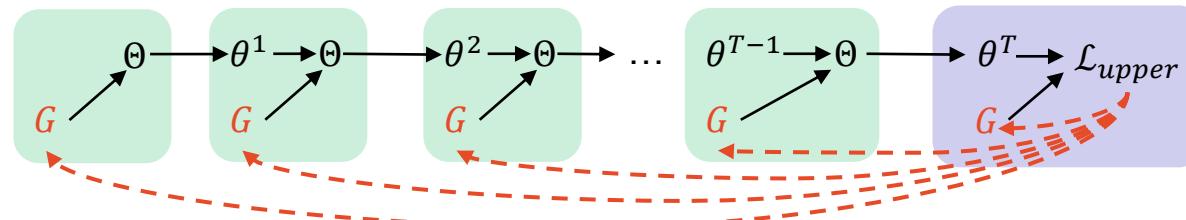
Performance of the trained classifier (i.e.,  $f(\theta^*)$ ) on the validation set ( $\mathcal{V}$ )

$$s.t. \theta^* = \arg \min_{\theta} - \sum_{i \in \mathcal{T}} \sum_{j=1}^c y_{ij} \ln f(G, \theta)[i, j]$$

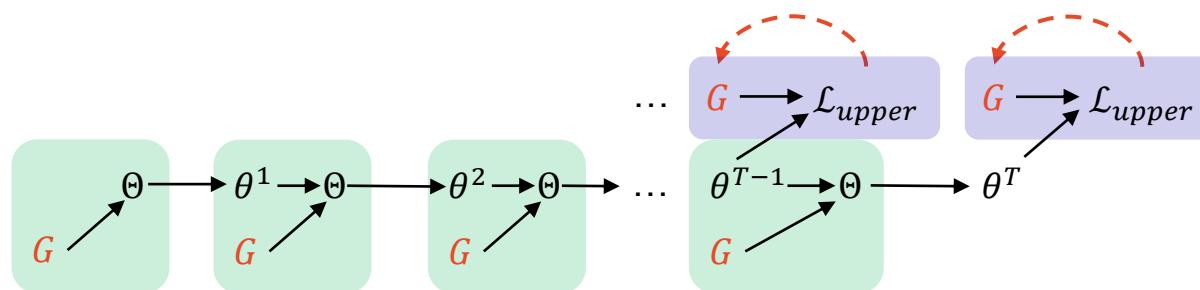
Performance of the classifier (i.e.,  $f(\theta)$ ) on the training set ( $\mathcal{T}$ )

# Hyper-gradient computation of GASOLINE

Iterative differentiation [1]: view the updating of lower-level problem as a **dynamic system** in  $T$  iterations (i.e.,  $\theta^* = \theta^T$ ):



First order approximation [2]:



$$\begin{aligned}\tilde{G} &= \arg \min_G \mathcal{L}_{upper}(G, \theta^*, \mathcal{V}, \mathbf{y}_{valid}) \\ \text{s.t. } \theta^* &= \arg \min_{\theta} \mathcal{L}_{lower}(G, \theta, \mathcal{T}, \mathbf{y}_{train})\end{aligned}$$

$\theta$ : iterative optimizer over  $\mathcal{L}_{lower}$  (e.g., GD)

—→ : gradient for updating  $G$

Sum all the gradient (i.e., —→) together as the hyper-gradient  $\nabla_G$  for updating  $G$

# Hyper-gradient-guided modification of GASOLINE

- Update  $G$  based on  $\nabla_G = \{\nabla_A, \nabla_X\}$  (update  $A$  as an example)

- (continuous update)** gradient descent:

- $$A \leftarrow A - \underbrace{\frac{B}{\sum_{i,j} |\nabla_A|_{[i,j]}}}_{\text{Budget constrained learning rate}} \cdot \nabla_A. \text{ Or}$$

Budget constrained learning rate

- (discrete update)** flip  $B$  entries in  $A$  whose indices obtain largest score in score matrix [1]:

- $$S = \underbrace{(-\nabla_A)}_{\text{preference}} \circ \underbrace{(1 - 2A)}_{\text{availability}}$$

preference availability

## Explanation of $S$

- Node  $i, j$  are connected,  $A[i, j] = 1$
- $(1 - 2A)[i, j] = -1$
- $-\nabla_A[i, j] > 0, S[i, j] < 0$
- $-\nabla_A[i, j] < 0, S[i, j] > 0$
- Dropping edge will get high score
- Reverse scenario for unconnected node pairs

# Low-rank speedup of GASOLINE (1)

- For updating  $\mathbf{A}$ ,  $\nabla_{\mathbf{A}}$  tends to be **dense**
- $\nabla_{\mathbf{A}} \mathcal{L}$ : time complexity  $O(n^2d)$ , space complexity  $O(n^2)$
- Recap the goal:  $\tilde{\mathbf{A}} = \mathbf{A} + \Delta\mathbf{A}$
- Key idea: the modification should be implemented over **a small set of nodes**, specifically:
  - decompose the incremental matrix:  $\Delta\mathbf{A} = \mathbf{U}\mathbf{V}'$

# Low-rank speedup of GASOLINE (2)

- Change optimization variable from  $\mathbf{A}$  to  $\mathbf{U}$  and  $\mathbf{V}$ . I.e., represent  $G = \{\mathbf{A} + \mathbf{UV}', \mathbf{X}\}$  in the formula:

$$\tilde{G} = \arg \min_G \mathcal{L}_{upper}(G, \theta^*, \mathcal{V}, \mathbf{y}_{valid})$$

s.t.  $\theta^*$

$$= \arg \min_{\theta} \mathcal{L}_{lower}(G, \theta, \mathcal{T}, \mathbf{y}_{train})$$

- Initialize  $\mathbf{U}$  and  $\mathbf{V}$  as zero-mean gaussian
- Compute  $\nabla_{\mathbf{U}}$  and  $\nabla_{\mathbf{V}}$  in the same paradigm as the full-version Gasoline
- Update  $\mathbf{U}$  and  $\mathbf{V}$  in a continuous way
- $\nabla_{\mathbf{U}}\mathcal{L}$  and  $\nabla_{\mathbf{V}}\mathcal{L}$ : time complexity  $O(nd^2 + md)$ , space complexity  $O(m + nd)$ 
  - $\nabla_{\mathbf{A}}\mathcal{L}$ : time complexity  $O(n^2d)$ , space complexity  $O(n^2)$

# Empirical performance of GASOLINE

Attack	Data	APPNP	GAT	Jaccard	SVD	RGCN	DE	LDS	G-DT	G-CF	G-DTCF
metattack	Cora	47.0±0.7	48.8±0.2	65.4±0.9	60.3±0.8	50.6±0.8	48.7±0.9	58.7±1.3	67.3±0.7	57.0±0.9	<b>68.8±0.9</b>
	Citeseer	49.4±2.2	62.4±0.7	57.1±1.0	49.5±0.8	55.5±1.4	50.1±2.3	58.2±2.3	<b>63.5±1.5</b>	58.4±1.5	62.2±1.0
	Polblogs	58.4±3.6	48.2±6.6	N/A	<b>79.1±2.4</b>	50.8±0.9	56.4±6.3	63.7±5.7	65.0±0.7	55.0±4.1	64.7±1.4
Nettack	Cora	60.7±1.2	54.2±2.3	63.7±1.4	52.9±2.8	56.5±1.1	60.8±1.0	64.5±2.4	64.5±2.2	63.9±2.4	<b>66.1±1.9</b>
	Citeseer	68.3±6.8	61.9±4.4	72.5±3.3	50.2±6.6	56.4±1.5	63.3±4.7	71.0±3.3	71.6±3.9	69.4±4.8	<b>74.3±1.6</b>
	Polblogs	90.5±1.0	91.1±0.7	N/A	<b>93.6±1.2</b>	93.1±0.2	89.1±2.4	91.1±1.8	92.3±1.6	90.3±0.7	92.4±1.7
random attack	Cora	74.3±0.4	58.1±1.0	75.1±0.5	72.6±0.3	68.9±0.4	73.9±0.6	76.6±0.4	77.1±0.3	<b>78.3±0.5</b>	77.8±0.2
	Citeseer	69.8±0.6	60.8±1.6	69.7±0.5	66.7±0.4	65.7±0.2	69.4±0.5	72.3±0.4	<b>73.8±0.2</b>	72.3±0.4	73.4±0.5
	Polblogs	74.7±2.8	<b>84.5±1.0</b>	N/A	83.3±2.8	81.7±0.9	75.9±1.4	73.2±2.8	73.4±4.1	77.1±1.6	77.6±2.9

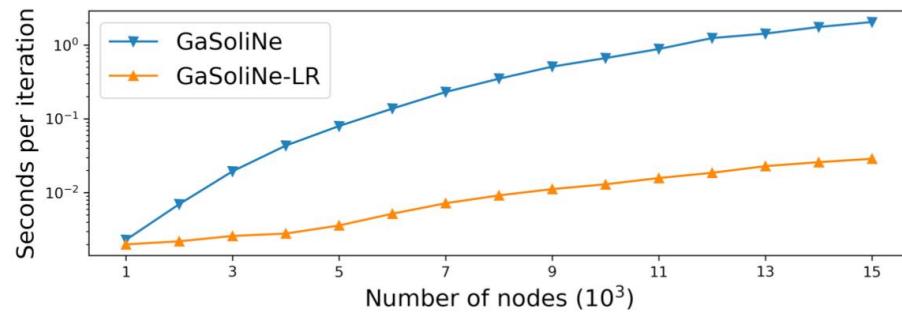
- Variants of Gasoline (**short as G**) perform best on Cora and CiteSeer
- Incorporating with **established methods**
- More experiments in our paper

Data	With GASOLINE?	GAT	SVD	RGCN
Cora	N	48.8±0.2	60.3±0.8	50.6±0.8
	Y	63.7±0.6	79.7±0.6	62.6±0.6
Citeseer	N	62.4±0.7	49.5±0.8	55.5±1.4
	Y	69.7±0.2	76.5±0.6	66.1±0.8
Polblogs	N	48.2±6.6	79.1±2.4	50.8±0.9
	Y	70.8±0.6	89.2±0.7	67.7±0.3

# Efficacy of the low-rank Gasoline

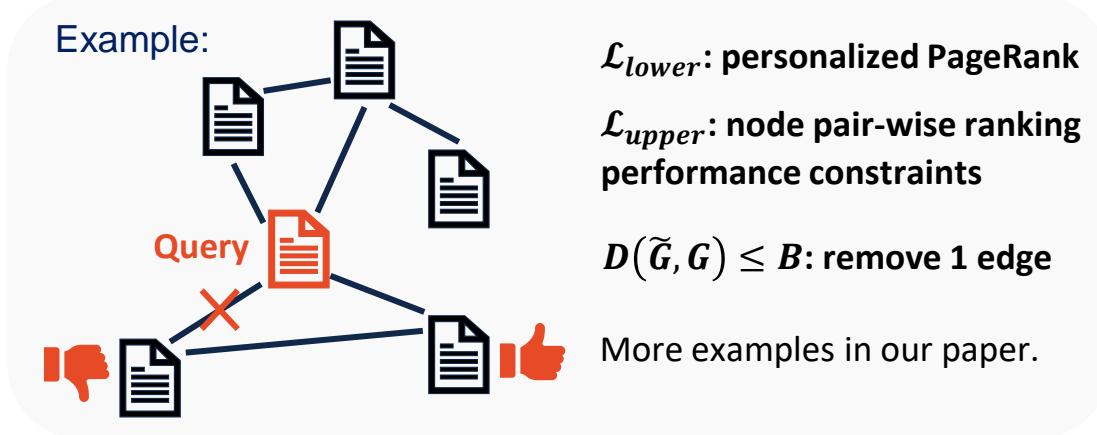
Data	Ptb Rate (%)	APPNP	GASOLINE	GASOLINE-LR
Cora	0	84.0±0.4	85.2±0.2	84.4±0.3
	5	74.1±0.7	77.4±0.5	75.0±0.3
	10	65.2±0.4	70.8±0.5	67.9±0.9
	15	58.2±1.1	67.1±0.8	65.3±0.8
	20	51.7±0.7	62.5±0.5	60.2±1.2
	25	47.0±0.7	57.3±0.6	57.1±0.5
Citeseer	0	71.8±0.4	74.7±0.2	73.4±0.2
	5	67.6±0.9	69.6±0.7	68.2±0.8
	10	61.8±0.8	66.3±1.0	63.9±0.4
	15	54.1±0.8	59.3±1.1	56.8±1.1
	20	51.0±1.2	56.5±0.9	55.3±0.9
	25	49.4±2.2	57.7±1.8	56.5±0.8
Polblogs	0	94.1±0.6	95.3±0.6	95.7±0.3
	5	70.1±0.6	73.8±0.9	93.4±0.3
	10	69.8±0.8	72.8±0.4	90.8±0.2
	15	67.5±0.5	70.1±1.2	88.7±0.3
	20	64.1±0.9	68.5±1.0	88.0±0.3
	25	57.0±3.6	64.8±2.1	89.9±0.5

- Comparable performance with Gasoline
- Strong performance on Polblogs
- Much more efficient compared with Gasoline



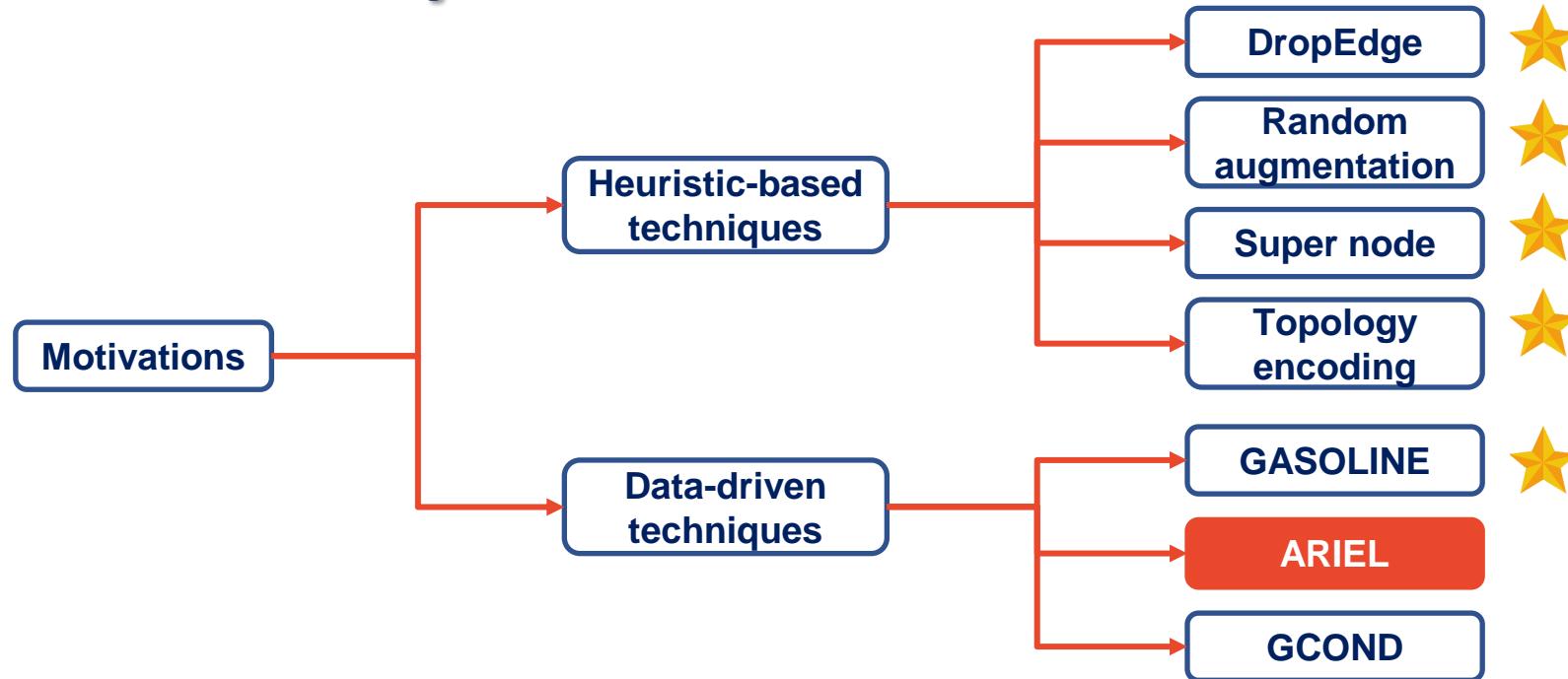
# Related works of GASOLINE

- GASOLINE [1] can work beyond node classification,  
E.g., supervised random walk [2,3]



- GASOLINE learns the adjacency matrix determinedly, and LDS [4] learns the graph generation matrix instead.

# Artificial dynamics in GNNs

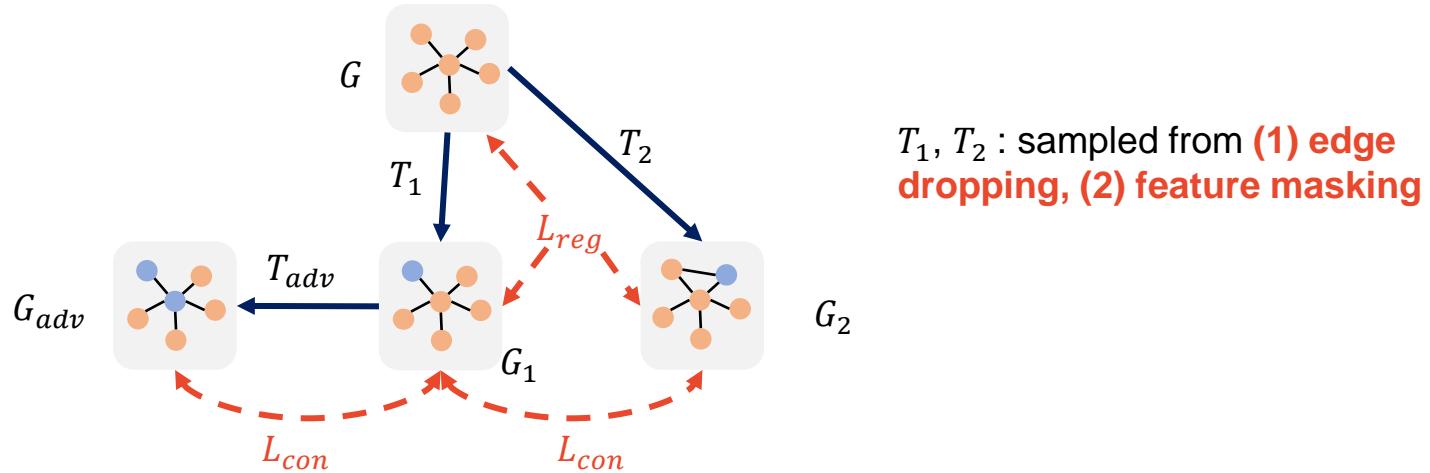


# Adversarial graph contrastive learning with information regularization (ARIEL)

**Task:** self-supervised node representation learning

**Artificial dynamic:** using adversarial attacks to generate an extra view of the given graph

**Goal:** improve graph contrastive learning performance



# Attack details of the ARIEL

$$G_{adv} = \arg \max_{G'} L_{con}(G_1, G')$$

s. t.  $\sum_{i,j} |\mathbf{A}'[i,j] - \mathbf{A}[i,j]| \leq \Delta_A, \sum_{i,j} |\mathbf{X}'[i,j] - \mathbf{X}[i,j]| \leq \Delta_X$

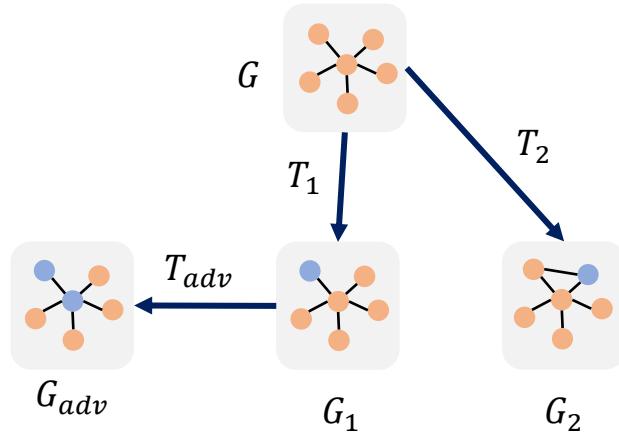
$\mathbf{L}_A$ : perturbations on edges;  $\mathbf{L}_X$ : perturbations on features

$$\mathbf{A}_{adv} = \mathbf{A} + (\mathbf{1}_{n \times n} - \mathbf{I}_n - 2\mathbf{A}) \circ \mathbf{L}_A$$

$$\mathbf{X}_{adv} = \mathbf{X} + \mathbf{L}_X$$

$\mathbf{L}_A \in \{0,1\}^{n \times n}$  is symmetric and discrete,

its convex hull  $\tilde{\mathbf{L}}_A \in [0,1]^{n \times n}$



# Project gradient attack (PGD) details (2)

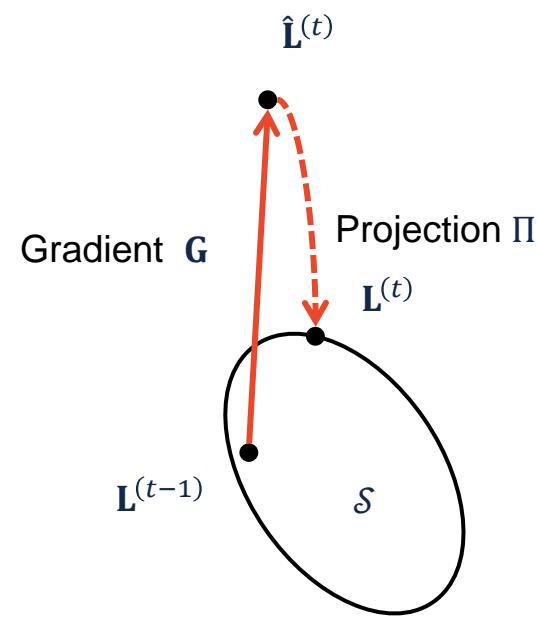
PGD attack updates:

$$\tilde{\mathbf{L}}_A^{(t)} = \Pi_{\mathcal{S}_A}(\tilde{\mathbf{L}}_A^{(t-1)} + \alpha \mathbf{G}_A^{(t)})$$

$$\mathbf{L}_X^{(t)} = \Pi_{\mathcal{S}_X}(\mathbf{L}_X^{(t-1)} + \beta \operatorname{sgn}(\mathbf{G}_X^{(t)}))$$

$\mathbf{G}_A^{(t)}$  and  $\mathbf{G}_X^{(t)}$ : gradients of **contrastive loss** with respect to

$\tilde{\mathbf{L}}_A^{(t-1)}$  and  $\mathbf{L}_X^{(t-1)}$



# Project gradient attack (PGD) details (2)

PGD attack updates:

$$\tilde{\mathbf{L}}_A^{(t)} = \Pi_{\mathcal{S}_A} \left( \tilde{\mathbf{L}}_A^{(t-1)} + \alpha \mathbf{G}_A^{(t)} \right), \mathbf{L}_X^{(t)} = \Pi_{\mathcal{S}_X} (\mathbf{L}_X^{(t-1)} + \beta \operatorname{sgn}(\mathbf{G}_X^{(t)}))$$

- $\Pi_{\mathcal{S}_A}(\mathbf{Z}) = \begin{cases} P(\mathbf{Z} - \mu \mathbf{1}_{n \times n}), & \sum_{i,j} P(\mathbf{Z} - \mu \mathbf{1}_{n \times n})[i,j] = \Delta_A \text{ for } \mu > 0 \\ P(\mathbf{Z}), & \sum_{i,j} P(\mathbf{Z})[i,j] \leq \Delta_A \end{cases}$

( $P(\mathbf{Z})$  elementwisely clips  $\mathbf{Z}$  into  $[0,1]$ )

- $\Pi_{\mathcal{S}_X}$  clips  $\mathbf{L}_X$  into the range  $[-\delta_X, \delta_X]$  elementwisely

# Regularization of the ARIEL

Information regularization: to constrain the views from the random augmentation

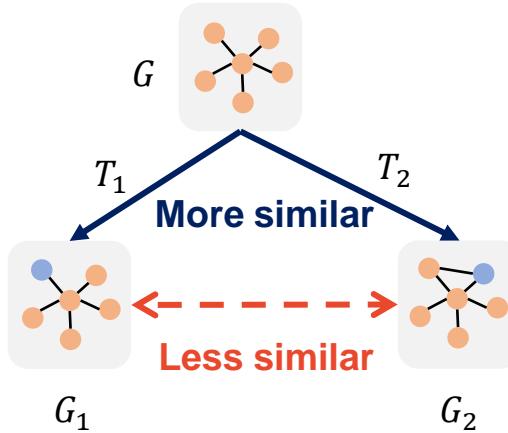
$$I_{i,1} = \text{similarity}(\mathbf{H}_1[i, :], \mathbf{H}[i, :])$$

$$I_{i,2} = \text{similarity}(\mathbf{H}_2[i, :], \mathbf{H}[i, :])$$

$$I_{i,12} = \text{similarity}(\mathbf{H}_1[i, :], \mathbf{H}_2[i, :])$$

$$d_i = 2I_{i,12} - (I_{i,1} + I_{i,2})$$

$$L_{reg}(G_1, G_2, G) = \frac{1}{n} \sum_{i=1}^n \max(d_i, 0)$$

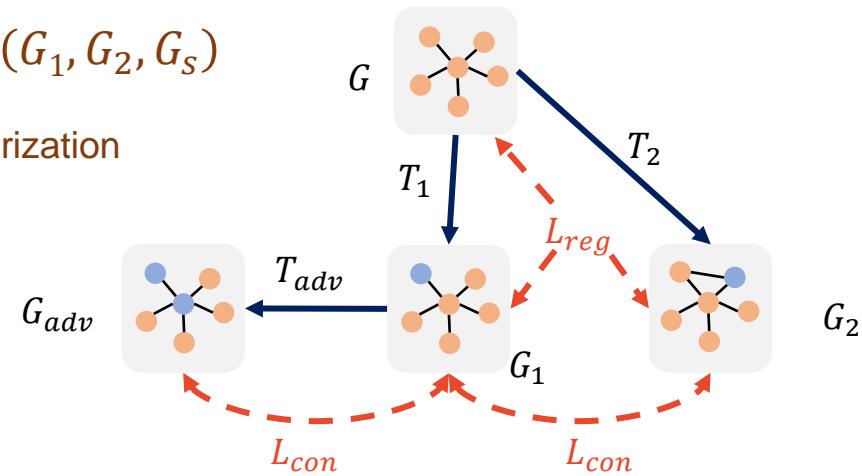


# Objective function of the ARIEL

$$L = L_{con}(G_1, G_2) + \varepsilon_1 L_{con}(G_1, G_{adv}) + \varepsilon_2 L_{reg}(G_1, G_2, G_s)$$

Contrastive loss Adversarial contrastive loss regularization

$$\text{s. t. } G_{adv} = \arg \max_{G'} L_{con}(G_1, G') , \\ \sum_{i,j} |\mathbf{A}'[i,j] - \mathbf{A}[i,j]| \leq \Delta_{\mathbf{A}}, \\ \sum_{i,j} |\mathbf{X}'[i,j] - \mathbf{X}[i,j]| \leq \Delta_{\mathbf{X}}$$



# Empirical performance (node classification accuracy) of the ARIEL

Method	Cora	CiteSeer	Amazon-Computers	Amazon-Photo	Coauthor-CS	Coauthor-Physics
GCN	$84.14 \pm 0.68$	$69.02 \pm 0.94$	$88.03 \pm 1.41$	$92.65 \pm 0.71$	$92.77 \pm 0.19$	$95.76 \pm 0.11$
GAT	$83.18 \pm 1.17$	$69.48 \pm 1.04$	$85.52 \pm 2.05$	$91.35 \pm 1.70$	$90.47 \pm 0.35$	$94.82 \pm 0.21$
DeepWalk+features	$79.82 \pm 0.85$	$67.14 \pm 0.81$	$86.23 \pm 0.37$	$90.45 \pm 0.45$	$85.02 \pm 0.44$	$94.57 \pm 0.20$
DGI	$84.24 \pm 0.75$	$69.12 \pm 1.29$	$88.78 \pm 0.43$	$92.57 \pm 0.23$	$92.26 \pm 0.12$	$95.38 \pm 0.07$
GMI	$82.43 \pm 0.90$	$70.14 \pm 1.00$	$83.57 \pm 0.40$	$88.04 \pm 0.59$	OOM	OOM
MVGRL	<b><math>84.39 \pm 0.77</math></b>	$69.85 \pm 1.54$	$89.02 \pm 0.21$	$92.92 \pm 0.25$	$92.22 \pm 0.22$	$95.49 \pm 0.17$
GRACE	$83.40 \pm 1.08$	$69.47 \pm 1.36$	$87.77 \pm 0.34$	$92.62 \pm 0.25$	$93.06 \pm 0.08$	$95.64 \pm 0.08$
GCA-DE	$82.57 \pm 0.87$	$72.11 \pm 0.98$	$88.10 \pm 0.33$	$92.87 \pm 0.27$	$93.08 \pm 0.18$	$95.62 \pm 0.13$
GCA-PR	$82.54 \pm 0.87$	$72.16 \pm 0.73$	$88.18 \pm 0.39$	$92.85 \pm 0.34$	$93.09 \pm 0.15$	$95.58 \pm 0.12$
GCA-EV	$81.80 \pm 0.92$	$67.07 \pm 0.79$	$87.95 \pm 0.43$	$92.63 \pm 0.33$	$93.06 \pm 0.14$	$95.64 \pm 0.08$
<b>ARIEL</b>	$84.28 \pm 0.96$	<b><math>72.74 \pm 1.10</math></b>	<b><math>91.13 \pm 0.34</math></b>	<b><math>94.01 \pm 0.23</math></b>	<b><math>93.83 \pm 0.14</math></b>	<b><math>95.98 \pm 0.05</math></b>

- More consistent improvement with adversarial training
- Better augmentation qualities than heuristic-based methods

# Robustness of the ARIEL

Topology attack with Metattack [1] and random mask attributes

Method	Cora	CiteSeer	Amazon-Computers	Amazon-Photos	Coauthor-CS	Coauthor-Physics
GCN	$80.03 \pm 0.91$	$62.98 \pm 1.20$	$84.10 \pm 1.05$	$91.72 \pm 0.94$	$80.32 \pm 0.59$	$87.47 \pm 0.38$
GAT	$79.49 \pm 1.29$	$63.30 \pm 1.11$	$81.60 \pm 1.59$	$90.66 \pm 1.62$	$77.75 \pm 0.80$	$86.65 \pm 0.41$
DeepWalk+features	$74.12 \pm 1.02$	$63.20 \pm 0.80$	$79.08 \pm 0.67$	$88.06 \pm 0.41$	$49.30 \pm 1.23$	$79.26 \pm 1.38$
DGI	$80.84 \pm 0.82$	$64.25 \pm 0.96$	$83.36 \pm 0.55$	$91.27 \pm 0.29$	$78.73 \pm 0.50$	$85.88 \pm 0.37$
GMI	$79.17 \pm 0.76$	$65.37 \pm 1.03$	$77.42 \pm 0.59$	$89.44 \pm 0.47$	$80.92 \pm 0.64$	$87.72 \pm 0.45$
MVGRL	<b><math>80.90 \pm 0.75</math></b>	$64.81 \pm 1.53$	$83.76 \pm 0.69$	$91.76 \pm 0.44$	$79.49 \pm 0.75$	$86.98 \pm 0.61$
GRACE	$78.55 \pm 0.81$	$63.17 \pm 1.81$	$84.74 \pm 1.13$	$91.26 \pm 0.37$	$80.61 \pm 0.63$	$85.71 \pm 0.38$
GCA	$76.79 \pm 0.97$	$64.89 \pm 1.33$	$85.05 \pm 0.51$	$91.71 \pm 0.34$	$82.72 \pm 0.58$	$89.00 \pm 0.31$
<b>ARIEL</b>	$80.33 \pm 1.25$	<b><math>69.13 \pm 0.94</math></b>	<b><math>88.61 \pm 0.46</math></b>	<b><math>92.99 \pm 0.21</math></b>	<b><math>84.43 \pm 0.59</math></b>	<b><math>89.09 \pm 0.31</math></b>

ARIEL is consistently robust on attacked graphs



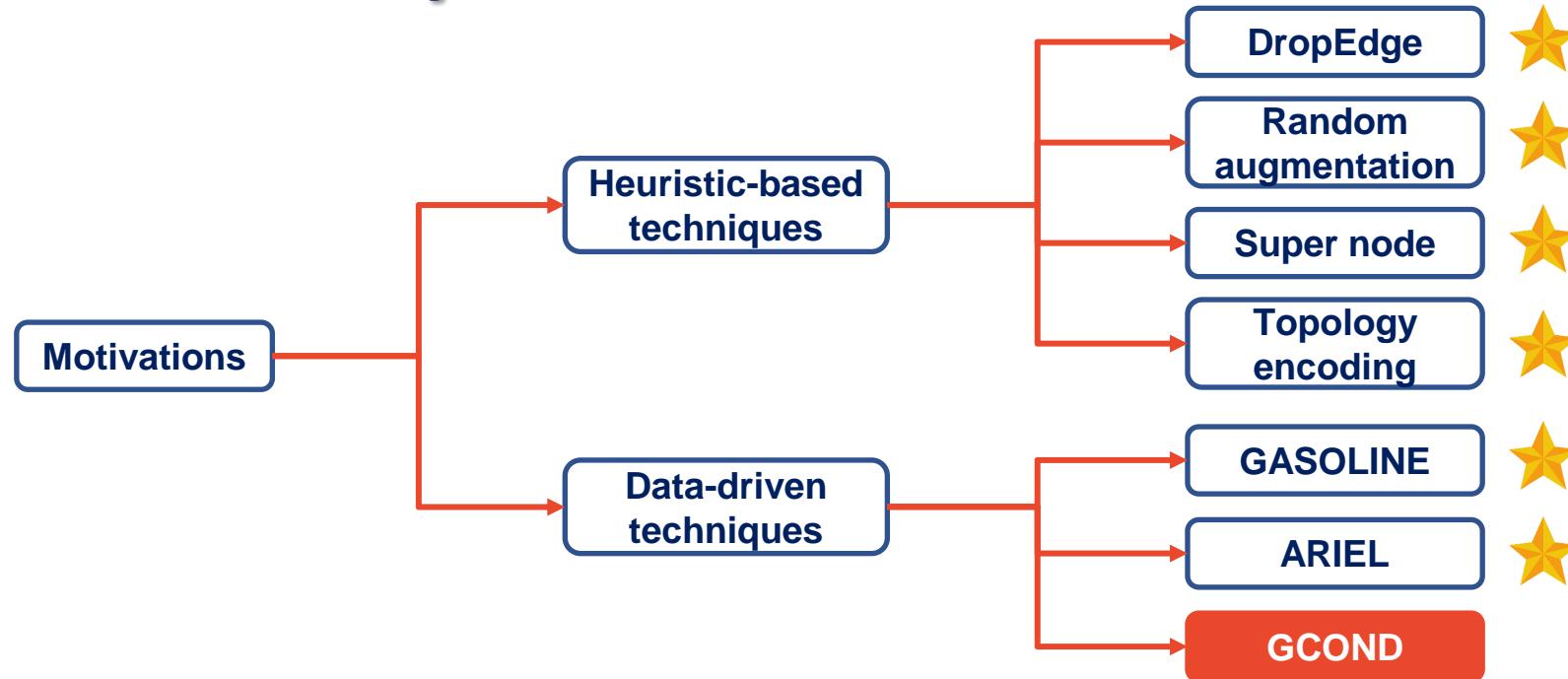
[1] Zügner et al. Adversarial Attacks on Graph Neural Networks via Meta Learning. ICLR 2019.

[2] Feng, Shengyu, et al. "Adversarial graph contrastive learning with information regularization." Proceedings of the ACM Web Conference 2022. 2022.

# Related works of the ARIEL

- Graph (virtual) **adversarial training**, e.g., [1,2]
- Adversarial graph augmentation to improve graph contrastive learning [3]: on **graph-level representation task**.
- Augmenting **hyper-graph** for contrastive learning in a data-driven fashion [4]

# Artificial dynamics in GNNs



# Graph condensation for graph neural networks (GCOND)

**Task:** node classification

**Artificial dynamic:** synthesizing a small graph on which GNN's **training gradient matches** the gradient on the original graph

**Goal:** Save storage space and training time for graph learning

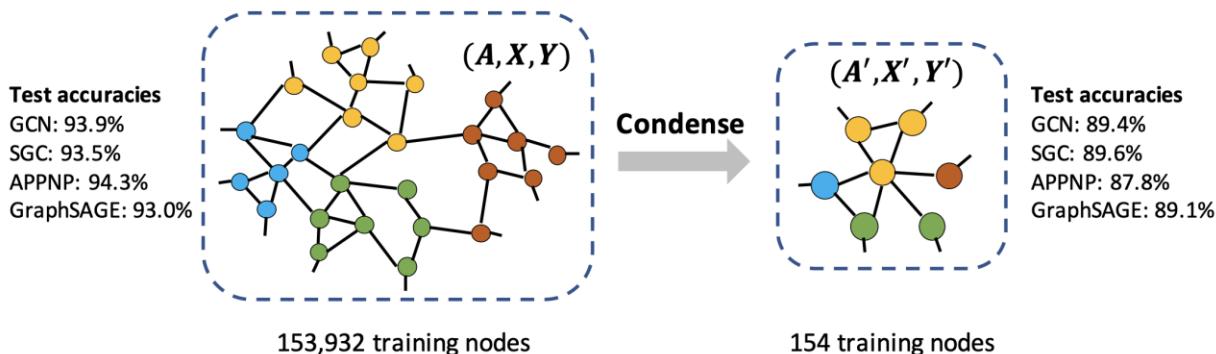


Figure credit: [1]

# Objective function of the GCOND

Given a graph  $\mathcal{T} = \{\mathbf{A}, \mathbf{X}, \mathbf{Y}\}$ , find the optimized synthetic graph  $\mathcal{S} = \{\mathbf{A}', \mathbf{X}', \mathbf{Y}'\}$ :

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[ \sum_t Distance(\nabla_{\theta} \mathcal{L}(GNN_{\theta_t}(\mathbf{A}', \mathbf{X}'), \mathbf{Y}'), \nabla_{\theta} \mathcal{L}(GNN_{\theta_t}(\mathbf{A}, \mathbf{X}), \mathbf{Y})) \right]$$

The choice of Distance function is flexible. [1] choices the **sum of the cosine distance** of every GNN layer.

# Optimization details of the GCOND

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[ \sum_t Distance(\nabla_{\theta} \mathcal{L}(GNN_{\theta_t}(\mathbf{A}', \mathbf{X}'), \mathbf{Y}'), \nabla_{\theta} \mathcal{L}(GNN_{\theta_t}(\mathbf{A}, \mathbf{X}), \mathbf{Y})) \right]$$

- Mini-batch: The given graph  $\mathcal{T} = \{\mathbf{A}, \mathbf{X}, \mathbf{Y}\}$  might be huge. In practice,  $\mathcal{T} = \{\mathbf{A}, \mathbf{X}, \mathbf{Y}\}$  can be replaced with a **sampled** graph  $\mathcal{T}_s = \{\mathbf{A}_s, \mathbf{X}_s, \mathbf{Y}_s\}$
- The search space of  $\mathbf{A}'$  is **quadratic** w.r.t. the number of nodes. To decrease the optimization variable, [1] proposes to encode  $\mathbf{A}'$  as  $MLP_{\phi}(\mathbf{X}')$ :

$$\mathbf{A}'[i, j] = \text{Sigmoid}\left(\frac{MLP_{\phi}(\mathbf{X}'[i: ,], \mathbf{X}'[j: ]) + MLP_{\phi}(\mathbf{X}'[j: ,], \mathbf{X}'[i: ])}{2}\right)$$

← **Node i/j's features**

# A graph-less version: GCOND-X

$$\min_{\mathcal{S}} \mathbb{E}_{\theta_0 \sim P_{\theta_0}} \left[ \sum_t Distance(\nabla_{\theta} \mathcal{L}(GNN_{\theta_t}(\mathbf{A}', \mathbf{X}'), \mathbf{Y}'), \nabla_{\theta} \mathcal{L}(GNN_{\theta_t}(\mathbf{A}, \mathbf{X}), \mathbf{Y})) \right]$$

- Explicitly set  $\mathbf{A}'$  as  $\mathbf{I}$
- A “special” representation learning method. Not only the topology information is encoded, but the # of nodes is also condensed.

# Empirical performance of the GCOND

Table 2: GCOND and GCOND-X achieves promising performance in comparison to baselines even with extremely large reduction rates. We report transductive performance on Citeseer, Cora, Ogbn-arxiv; inductive performance on Flickr, Reddit. Performance is reported as test accuracy (%).

Dataset	Ratio ( $r$ )	Baselines					Proposed		
		Random ( $\mathbf{A}', \mathbf{X}'$ )	Herding ( $\mathbf{A}', \mathbf{X}'$ )	K-Center ( $\mathbf{A}', \mathbf{X}'$ )	Coarsening ( $\mathbf{A}', \mathbf{X}'$ )	DC-Graph ( $\mathbf{X}'$ )	GCOND-X ( $\mathbf{X}'$ )	GCOND ( $\mathbf{A}', \mathbf{X}'$ )	Whole Dataset
Citeseer	0.9%	54.4 $\pm$ 4.4	57.1 $\pm$ 1.5	52.4 $\pm$ 2.8	52.2 $\pm$ 0.4	66.8 $\pm$ 1.5	<b>71.4</b> $\pm$ 0.8	70.5 $\pm$ 1.2	
	1.8%	64.2 $\pm$ 1.7	66.7 $\pm$ 1.0	64.3 $\pm$ 1.0	59.0 $\pm$ 0.5	66.9 $\pm$ 0.9	69.8 $\pm$ 1.1	<b>70.6</b> $\pm$ 0.9	71.7 $\pm$ 0.1
	3.6%	69.1 $\pm$ 0.1	69.0 $\pm$ 0.1	69.1 $\pm$ 0.1	65.3 $\pm$ 0.5	66.3 $\pm$ 1.5	69.4 $\pm$ 1.4	<b>69.8</b> $\pm$ 1.4	
Cora	1.3%	63.6 $\pm$ 3.7	67.0 $\pm$ 1.3	64.0 $\pm$ 2.3	31.2 $\pm$ 0.2	67.3 $\pm$ 1.9	75.9 $\pm$ 1.2	<b>79.8</b> $\pm$ 1.3	
	2.6%	72.8 $\pm$ 1.1	73.4 $\pm$ 1.0	73.2 $\pm$ 1.2	65.2 $\pm$ 0.6	67.6 $\pm$ 3.5	75.7 $\pm$ 0.9	<b>80.1</b> $\pm$ 0.6	81.2 $\pm$ 0.2
	5.2%	76.8 $\pm$ 0.1	76.8 $\pm$ 0.1	76.7 $\pm$ 0.1	70.6 $\pm$ 0.1	67.7 $\pm$ 2.2	76.0 $\pm$ 0.9	<b>79.3</b> $\pm$ 0.3	
Ogbn-arxiv	0.05%	47.1 $\pm$ 3.9	52.4 $\pm$ 1.8	47.2 $\pm$ 3.0	35.4 $\pm$ 0.3	58.6 $\pm$ 0.4	<b>61.3</b> $\pm$ 0.5	59.2 $\pm$ 1.1	
	0.25%	57.3 $\pm$ 1.1	58.6 $\pm$ 1.2	56.8 $\pm$ 0.8	43.5 $\pm$ 0.2	59.9 $\pm$ 0.3	<b>64.2</b> $\pm$ 0.4	63.2 $\pm$ 0.3	71.4 $\pm$ 0.1
	0.5%	60.0 $\pm$ 0.9	60.4 $\pm$ 0.8	60.3 $\pm$ 0.4	50.4 $\pm$ 0.1	59.5 $\pm$ 0.3	63.1 $\pm$ 0.5	<b>64.0</b> $\pm$ 0.4	
Flickr	0.1%	41.8 $\pm$ 2.0	42.5 $\pm$ 1.8	42.0 $\pm$ 0.7	41.9 $\pm$ 0.2	46.3 $\pm$ 0.2	45.9 $\pm$ 0.1	<b>46.5</b> $\pm$ 0.4	
	0.5%	44.0 $\pm$ 0.4	43.9 $\pm$ 0.9	43.2 $\pm$ 0.1	44.5 $\pm$ 0.1	45.9 $\pm$ 0.1	45.0 $\pm$ 0.2	<b>47.1</b> $\pm$ 0.1	47.2 $\pm$ 0.1
	1%	44.6 $\pm$ 0.2	44.4 $\pm$ 0.6	44.1 $\pm$ 0.4	44.6 $\pm$ 0.1	45.8 $\pm$ 0.1	45.0 $\pm$ 0.1	<b>47.1</b> $\pm$ 0.1	
Reddit	0.05%	46.1 $\pm$ 4.4	53.1 $\pm$ 2.5	46.6 $\pm$ 2.3	40.9 $\pm$ 0.5	88.2 $\pm$ 0.2	<b>88.4</b> $\pm$ 0.4	88.0 $\pm$ 1.8	
	0.1%	58.0 $\pm$ 2.2	62.7 $\pm$ 1.0	53.0 $\pm$ 3.3	42.8 $\pm$ 0.8	89.5 $\pm$ 0.1	89.3 $\pm$ 0.1	<b>89.6</b> $\pm$ 0.7	93.9 $\pm$ 0.0
	0.2%	66.3 $\pm$ 1.9	71.0 $\pm$ 1.6	58.5 $\pm$ 2.1	47.4 $\pm$ 0.9	<b>90.5</b> $\pm$ 1.2	88.8 $\pm$ 0.4	90.1 $\pm$ 0.5	

- Empirically, GCOND is the state-of-the-art method for condensing a given graph.
- Even the structure-free version GCOND-X can get decent performance.

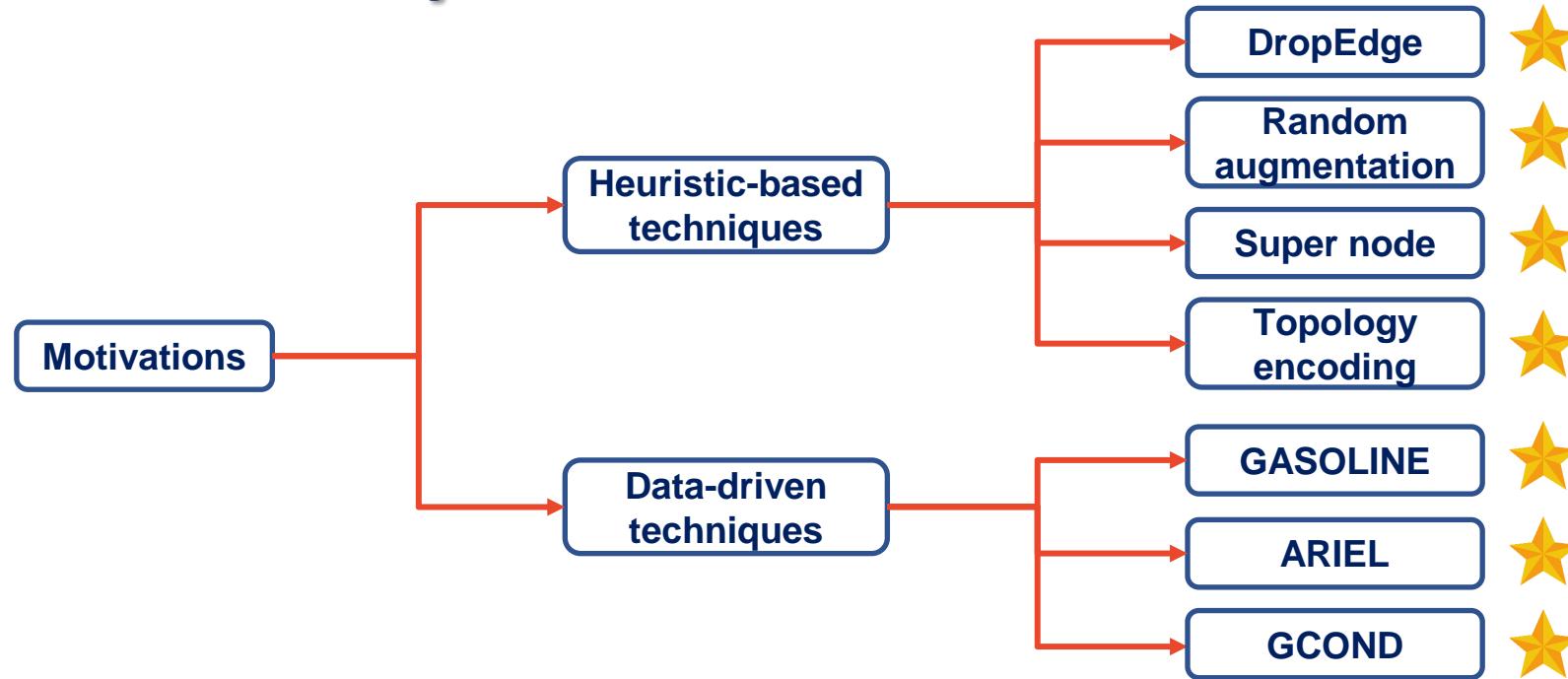
# Applicability of the GCOND

Table 3: Graph condensation can work well with different architectures. Avg. stands for the average test accuracy of APPNP, Cheby, GCN, GraphSAGE and SGC. SAGE stands for GraphSAGE.

	Methods	Data	MLP	GAT	APPNP	Cheby	GCN	SAGE	SGC	Avg.
Citeseer $r = 1.8\%$	DC-Graph	$\mathbf{X}'$	66.2	-	66.4	64.9	66.2	65.9	69.6	66.6
	GCOND-X	$\mathbf{X}'$	69.6	-	69.7	70.6	69.7	69.2	71.6	70.2
	GCOND	$\mathbf{A}', \mathbf{X}'$	63.9	55.4	69.6	68.3	70.5	66.2	70.3	69.0
Cora $r = 2.6\%$	DC-Graph	$\mathbf{X}'$	67.2	-	67.1	67.7	67.9	66.2	72.8	68.3
	GCOND-X	$\mathbf{X}'$	76.0	-	77.0	74.1	75.3	76.0	76.1	75.7
	GCOND	$\mathbf{A}', \mathbf{X}'$	73.1	66.2	78.5	76.0	80.1	78.2	79.3	78.4
Ogbn-arxiv $r = 0.25\%$	DC-Graph	$\mathbf{X}'$	59.9	-	60.0	55.7	59.8	60.0	60.4	59.2
	GCOND-X	$\mathbf{X}'$	64.1	-	61.5	59.5	64.2	64.4	64.7	62.9
	GCOND	$\mathbf{A}', \mathbf{X}'$	62.2	60.0	63.4	54.9	63.2	62.6	63.7	61.6
Flickr $r = 0.5\%$	DC-Graph	$\mathbf{X}'$	43.1	-	45.7	43.8	45.9	45.8	45.6	45.4
	GCOND-X	$\mathbf{X}'$	42.1	-	44.6	42.3	45.0	44.7	44.4	44.2
	GCOND	$\mathbf{A}', \mathbf{X}'$	44.8	40.1	45.9	42.8	47.1	46.2	46.1	45.6
Reddit $r = 0.1\%$	DC-Graph	$\mathbf{X}'$	50.3	-	81.2	77.5	89.5	89.7	90.5	85.7
	GCOND-X	$\mathbf{X}'$	40.1	-	78.7	74.0	89.3	89.3	91.0	84.5
	GCOND	$\mathbf{A}', \mathbf{X}'$	42.5	60.2	87.8	75.5	89.4	89.1	89.6	86.3

- The condensed graph from GCOND and GCOND-X can be applied to various downstream models.

# Artificial dynamics in GNNs



# Contents

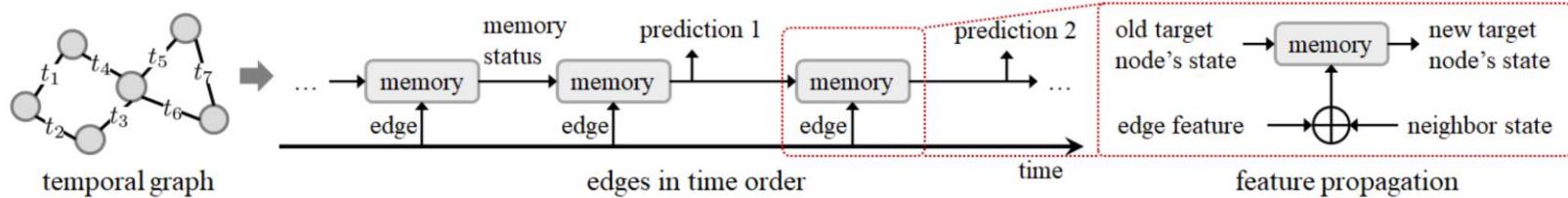
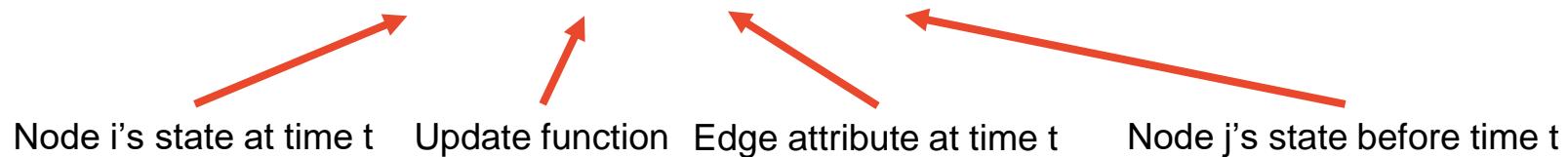
- Part I – Introduction 
- Part II – Natural Dynamics in GNNs 
- 30 mins Coffee Break 
- Part III – Artificial Dynamics in GNNs 
- Part IV – Challenges and Future Directions
- Q&A

# Open questions (1)

**Question 1:** how to augment the temporal graphs, i.e., the temporal or the topological features?

**Recap:** Encoding temporal graphs [1]

$$\mathbf{s}_i(t) = \text{mem}(\mathbf{s}_i(t^-), \{\mathbf{e}_{ij}(t) \oplus \mathbf{s}_j(t^-) \mid j \in N(i)\})$$



# Adaptive data augmentation on temporal graphs (MeTA)

**Task:** Edge prediction and node classification

**(Heuristic) artificial dynamics:** (1) perturbing time:  $t_{new} = t + \tau$ ,  $\tau \sim N(0, \sigma^2)$ , (2) edge dropping, (3) edge adding with noisy time

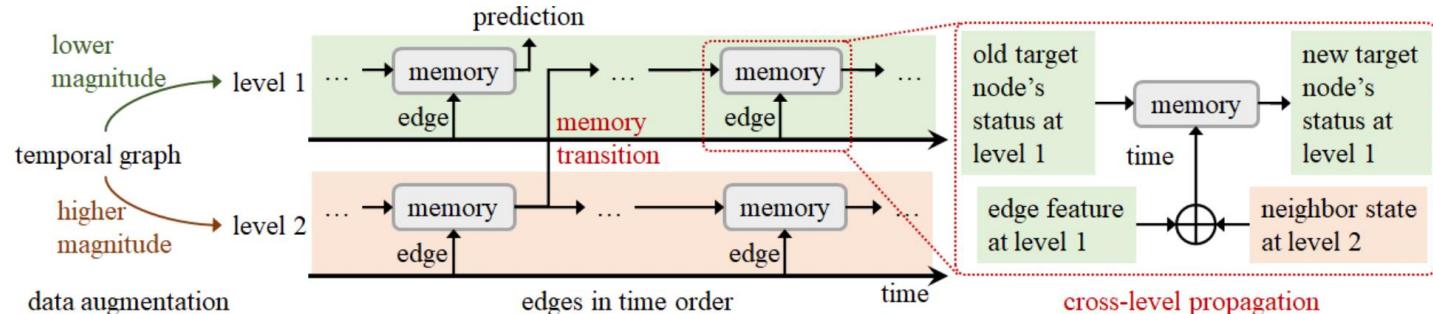
**Goal:** improving temporal GNNs' performance

**Key idea:** (1) augmenting both the **topological component** and **the time features**; (2) merging representations from graphs with **different augmentation magnitudes**

# Details of MeTA

E.g., perturb 10% edges (less random)

E.g., perturb 80% edges (more random)



$$\mathbf{s}_i^{(l)}(t) = \text{mem}(\mathbf{s}_i^{(l)}(t^-), \mathbf{e}_{ij}^{(l)}(t) \| \mathbf{s}_j^{\min(l+1,L)}(t^-)),$$

Node embedding from the less random graph

“adaptive”

Node embedding from the more random graph

# Applicability of MeTA

Table 2: Test accuracy and average precision (AP) of transductive edge prediction. We conduct 100 trials with random weight initialization. Mean (%) and standard deviations are reported. The best results in each column are highlighted in **bold** font.

Method	MOOC		Reddit		Wikipedia	
	Accuracy	AP	Accuracy	AP	Accuracy	AP
CTDNE [21]	65.34 ± 0.7	74.29 ± 0.6	73.76 ± 0.5	91.41 ± 0.3	79.42 ± 0.4	92.17 ± 0.5
JODIE [16]	76.45 ± 0.6	83.87 ± 0.4	90.91 ± 0.3	97.11 ± 0.3	87.04 ± 0.4	94.62 ± 0.5
TGAT [39]	75.20 ± 0.5	82.66 ± 0.4	92.92 ± 0.3	98.12 ± 0.2	88.14 ± 0.2	95.34 ± 0.1
DyRep [28]	73.36 ± 0.4	81.75 ± 0.3	92.11 ± 0.2	97.98 ± 0.1	87.77 ± 0.2	94.59 ± 0.2
TGN [25]	81.38 ± 0.6	89.79 ± 0.5	92.56 ± 0.2	98.70 ± 0.1	89.51 ± 0.4	98.46 ± 0.1
DyRep + MeTA (Ours)	76.21 ± 0.4	84.18 ± 0.3	93.04 ± 0.3	98.62 ± 0.1	88.92 ± 0.2	95.63 ± 0.2
TGN + MeTA (Ours)	<b>83.84 ± 0.5</b>	<b>92.03 ± 0.3</b>	<b>94.19 ± 0.2</b>	<b>99.08 ± 0.1</b>	<b>91.34 ± 0.3</b>	<b>98.87 ± 0.1</b>

- The method can work with **different backbone models** (e.g., DyRep and TGN)
- It can work on both the **edge-level tasks** and the **node-level tasks** (omitted in this slides for brevity)

# Open questions (2)

**Question 2:** how to augment the temporal graph in a data-driven way?

**Recap:** [1] perturb the edge time heuristically:  $t_{new} = t + \tau$ ,  $\tau \sim N(0, \sigma^2)$ .

**Challenge:** If we aim to augment (i.e., optimize) the edge time in a data-driven way (e.g., **gradient-based method**), a given optimization objective must be **differentiable** w.r.t. the “time” parameters.

$$\mathbf{s}_i(t) = \text{mem}(\mathbf{s}_i(t^-), \{\mathbf{e}_{ij}(t) \oplus \mathbf{s}_j(t^-) \mid j \in N(i)\})$$


Node i's state at time t    Update function    Edge attribute at time t    Node i's state before time t

# Open questions (3)

**Question 3:** How to ensure augmentation efficiency?

**Challenges:**

- Even if we do not augment the time parameter and view the temporal graph as a sequence of static graph snapshots, the method running time is linear w.r.t. # of time steps.
- The "optimal" augmented graph snapshots should be related. How can we incrementally get the optimal graph snapshots from previous time steps?

# Open questions (4)

**Question 4:** How to generate new graph snapshots given the previously observed ones?

**Challenges:**

- Highly open and there are many unknown hyper-parameters. E.g., How many graphs are generated? What is the size of generated graphs?
- How to ensure the performance of downstream tasks is improved given the generated graphs?

# Takeaways

- **Introduction**
- **Natural Dynamics in GNNs**
- **Artificial Dynamics in GNNs**
  - Heuristic-based techniques
  - Data-driven techniques
- **Open Questions and Challenges**
  - Augmentation of both the topological and temporal components
  - Augmenting the temporal modality in a data-driven way
  - How to maintain efficiency

# Resources

- **Surveys**

- Dongqi Fu and Jingrui He: Natural and Artificial Dynamics in Graphs: Concept, Progress, and Future. *Frontiers in Big Data* 2022
- Ding, Kaize, et al. "Data augmentation for deep graph learning: A survey." *ACM SIGKDD Explorations Newsletter* 24.2 (2022): 61-77.
- Zhao, Tong, et al. "Graph data augmentation for graph machine learning: A survey." *arXiv preprint arXiv:2202.08871* (2022).

- **Related Tutorials**

- Rozenshtein, Polina, and Aristides Gionis. "Mining temporal networks." *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*. 2019.
- Zhang, Chuxu, Jundong Li, and Meng Jiang. "Data Efficient Learning on Graphs." *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2021.
- Jin, Wei, et al. "Graph representation learning: foundations, methods, applications and systems." *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. 2021.

# Contents

- Part I – Introduction 
- Part II – Natural Dynamics in GNNs 
- 30 mins Coffee Break 
- Part III – Artificial Dynamics in GNNs 
- Part IV – Challenges and Future Directions 
- Q&A



# Thanks!



Dongqi Fu  
UIUC



Zhe Xu  
UIUC



Hanghang Tong  
UIUC



Jingrui He  
UIUC

{dongqif2, zhixu3, htong, jingrui}@illinois.edu

<https://github.com/DongqiFu/Natural-and-Artificial-Dynamics-in-GNNs-A-Tutorial>

