

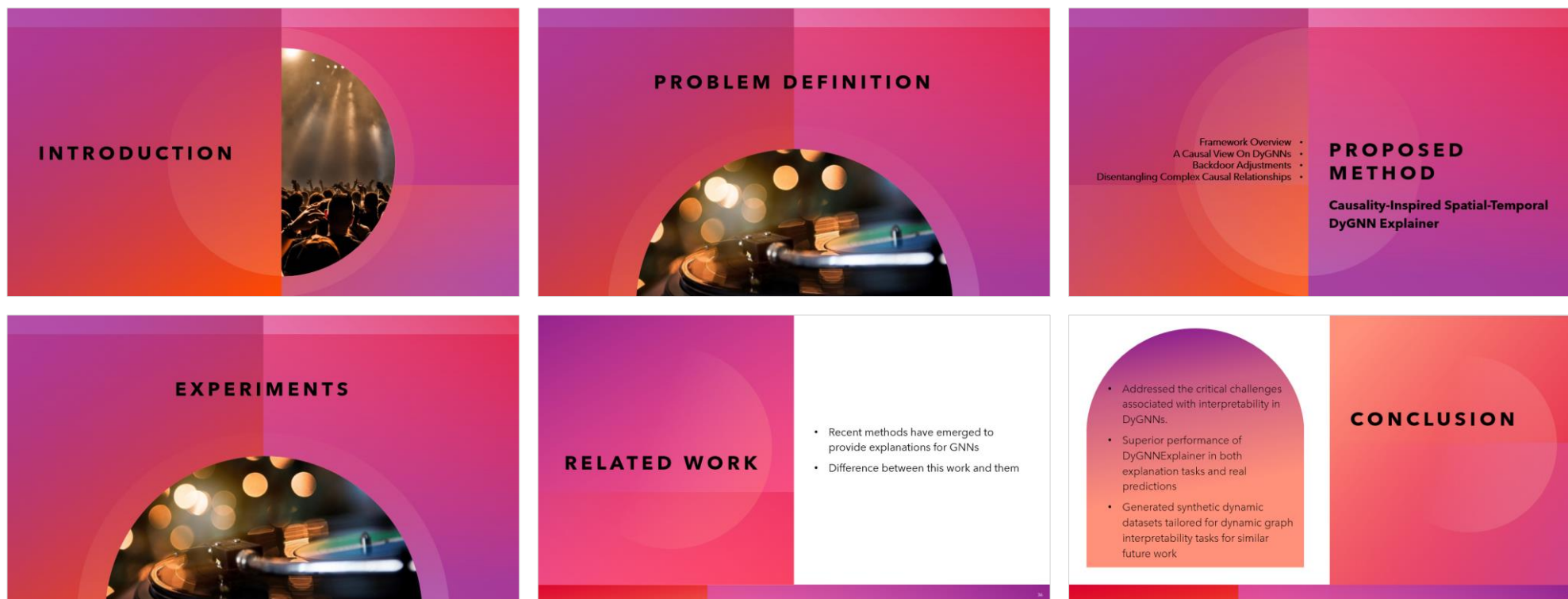
**CAUSALITY-INSPIRED  
SPATIAL-TEMPORAL EXPLANATIONS  
FOR DYNAMIC GRAPH NEURAL NETWORKS**

*Hooman Zolfaghari*

# AGENDA

- Introduction
- Problem Definition
- Proposed Method
- Experiments
- Related Work
- Conclusion

# ZOOM



# INTRODUCTION



# PROBLEM DEFINITION



# OVERVIEW OF THE TARGET MODEL

Dynamic Graph Neural Network (DyGNN)

- $f = f_d \circ f_a$
- $f_a: \mathcal{G}_{1:T} \rightarrow \mathcal{R}$ 
  - **Aggregation Function:** Captures temporal structures and feature patterns
  - **Input:** Dynamic graphs  $\mathcal{G}_{1:T}$  over  $T$  time steps.
  - **Output:** High-dimensional graph representation in  $\mathcal{R}$
- Downstream Task:  $f_d: \mathcal{R} \rightarrow \mathcal{Y}$ 
  - **Downstream Task Function:** Transforms graph representation to label space.
  - **Output:** Final label prediction in  $\mathcal{Y}$

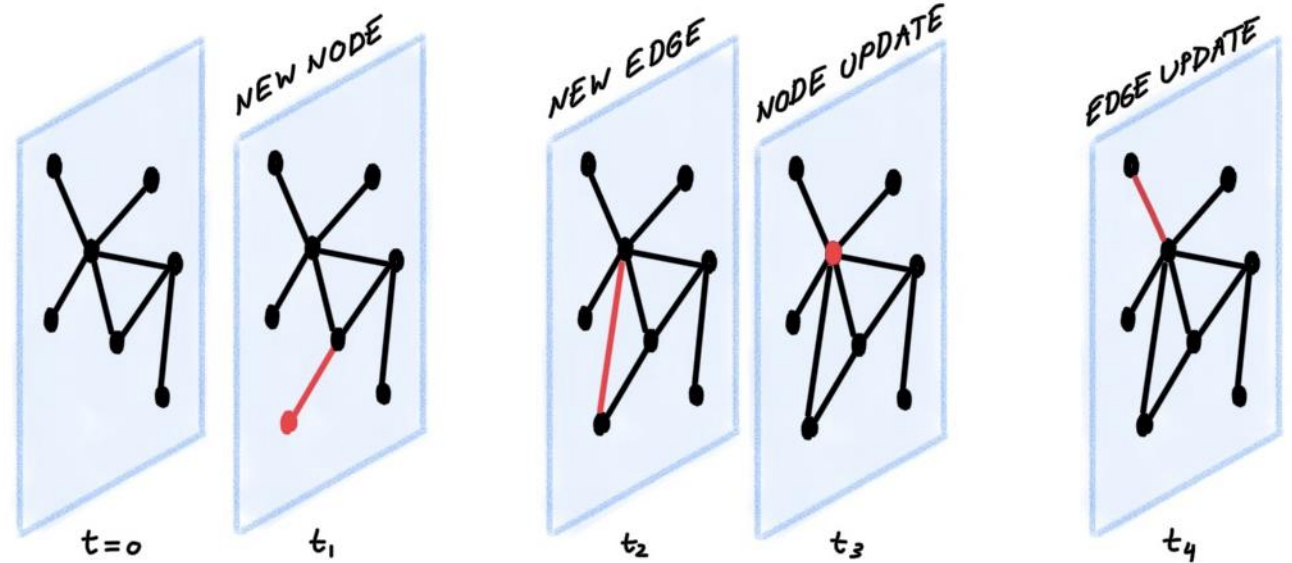


# OVERVIEW OF THE TARGET MODEL

## (CONT.)

### Dynamic Graph Structure

- Input at Each Time Step  $t$ :  $G_t = (X_t, A_t)$
- Node Attribute Matrix:  $X_t \in \mathbb{R}^{|V| \times D}$ 
  - $|V|$ : Number of nodes.
  - $D$ : Dimension of node attributes.
- Adjacency Matrix:  $A_t \in \mathbb{R}^{|V| \times |V|}$



# OBJECTIVE

Explanation Objective Critical Criteria for DyGNNs

## 1. Fidelity

- Explanations (dynamic subgraphs) should accurately reflect the model's behavior around the input graphs.
- Explanatory subgraphs, when fed back into the model, should produce similar predictions as the original dynamic graphs.

## 2. Interpretability

- Explanations should highlight the most important parts of the input while ignoring irrelevant components.
- **Spatial Interpretability:** Identifies key subgraphs within each graph
- **Temporal Interpretability:** Identifies key time steps contributing to the prediction



# OBJECTIVE (CONT.)

Model-Agnostic Explainer. Generative Model  $\mathcal{F}$

- Identifies aspects of the input that contribute to DyGNN's predictions.
- Meets both fidelity and interpretability criteria.
- **Model-Agnostic:** Can explain any black-box DyGNN without needing access to its internal mechanisms or ground-truth labels.
- **Focus:** providing spatial-temporal explanations (dynamic subgraph set) for dynamic graph structures.

- Framework Overview
- A Causal View On DyGNNs
- Backdoor Adjustments
- Disentangling Complex Causal Relationships

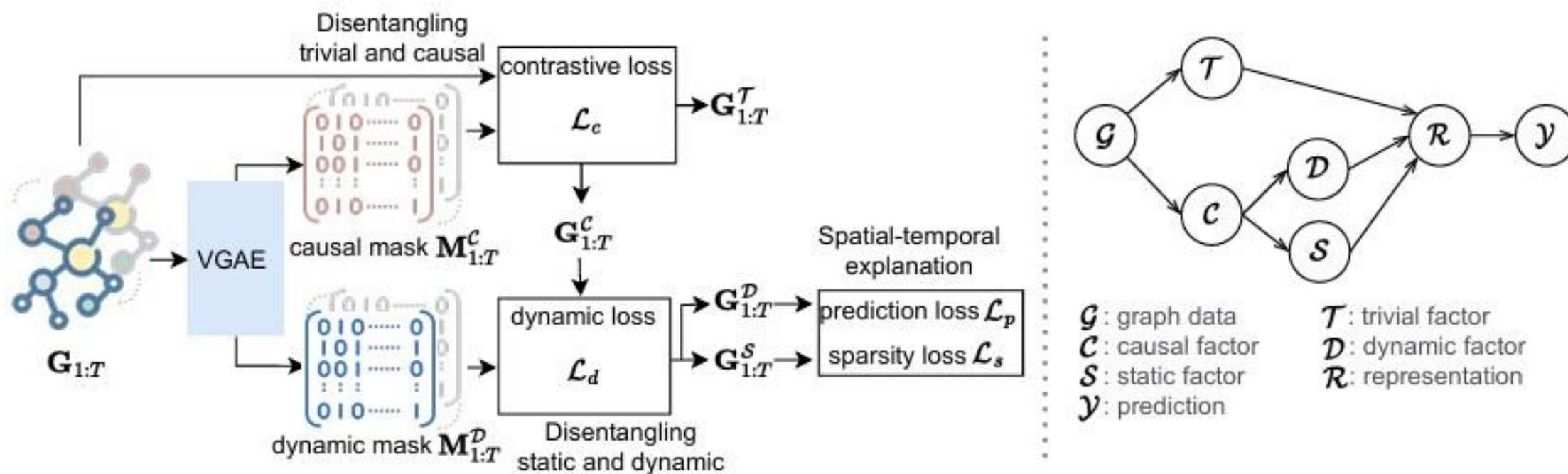
# PROPOSED METHOD

**Causality-Inspired Spatial-Temporal  
DyGNN Explainer**

# FRAMEWORK OVERVIEW

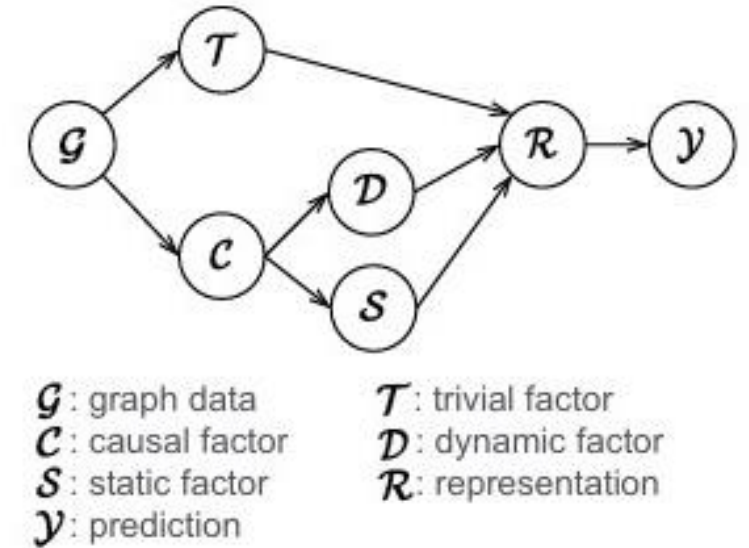
1. Construct a Structural Causal Model (SCM) , enabling a comprehensive understanding of dynamic graphs
2. generate causal and dynamic soft masks to enable **backdoor adjustment**, to **intervene** in the targeting causal and dynamic factors
3. Use:
  - Contrastive Loss: Separates trivial and causal relationships
  - Dynamic Loss: Disentangles static and dynamic relationships
4. Use Prediction & Sparsity Loss: Enhance prediction accuracy and interpretability

# FRAMEWORK OVERVIEW



# CAUSAL VIEW

- $\mathcal{T} \leftarrow \mathcal{G} \rightarrow \mathcal{C}$  :
  - $\mathcal{C}$ : Genuine causal relationships
  - $\mathcal{T}$ : Trivial relationships (data biases/spurious patterns)
- $\mathcal{T} \rightarrow \mathcal{R} \leftarrow \mathcal{C}$ 
  - $\mathcal{R}$ : high-dimensional representation of dynamic graph node data.
  - Uses both  $\mathcal{T}$  and  $\mathcal{C}$  to extract discriminative information.
- $\mathcal{D} \rightarrow \mathcal{R} \leftarrow \mathcal{S}$ : Causal relationships consist of dynamic relationship and static relationship
- $\mathcal{R} \rightarrow \mathcal{Y}$  : Ultimate aim of dynamic graph representation learning to predict graph properties



# CAUSAL VIEW

Backdoor Paths from the SCM:

- Path 1:  $\mathcal{C} \leftarrow \mathcal{G} \rightarrow \mathcal{T} \rightarrow \mathcal{R} \rightarrow \mathcal{Y}$ 
  - Issue:  $\mathcal{T}$  acts as a confounder between  $\mathcal{G}$  and  $\mathcal{Y}$
  - Consequence: Misleading correlation between  $\mathcal{C}$  and  $\mathcal{Y}$  leading to incorrect predictions
- Path 2:  $\mathcal{D} \leftarrow \mathcal{C} \rightarrow \mathcal{S} \rightarrow \mathcal{R} \rightarrow \mathcal{Y}$ 
  - Issue:  $\mathcal{S}$  acts as a confounder between  $\mathcal{C}$  and  $\mathcal{Y}$
- Thus:
  - It is important to block backdoor paths to ensure DyGNNs utilize genuine causal relationships without interference from confounders



# BACKDOOR ADJUSTMENT

- Safeguard DyGNNs against confounding factors and distinguish between dynamic and static relationships
- Focus: representation learning that eliminates the backdoor path
- Do-Calculus on  $\mathcal{T}$ :
  - Estimate  $P(\mathcal{Y}|\text{do}(\mathcal{C}))$  without interference from  $\mathcal{T}$
- Do-Calculus on  $\mathcal{D}$ :
  - Estimate  $P(\mathcal{Y}|\text{do}(\mathcal{D}))$  eliminating backdoor from  $\mathcal{S}$

# BACKDOOR ADJUSTMENT

Challenge:

- Can't directly employ the standard backdoor adjustment method due to confounder  $\mathcal{T}$
- Merge the estimation of  $P(\mathcal{Y}|do(\mathcal{C}))$  with that of  $P(\mathcal{Y}|do(\mathcal{D}))$  :
$$P(\mathcal{Y}|do(\mathcal{D})) = \sum P(\mathcal{Y}|do(\mathcal{C}))P(\mathcal{S}) = \sum P(\mathcal{S}) \sum P(\mathcal{Y}|\mathcal{G})P(\mathcal{T})$$
- No explicit information is available for the identification of trivial, dynamic, and static relationships.

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

- Casual soft masks for Casual relationship at  $t^{th}$  time step as  $M_t^c \in \mathbb{R}^{|V| \times |V|}$ .
- Each element represents an attention score typically in  $[0,1]$ .
- Complementary Mask:  $\bar{M} = \mathbf{1} - M$
- Partition dynamic graph set as:
  - Casual set  $G_{1:T}^c = (X_{1:T}, A_{1:T} \oplus M_{1:T}^c)$
  - Trivial set  $G_{1:T}^T = (X_{1:T}, A_{1:T} \oplus \bar{M}_{1:T}^c)$ 
    - Elementwise dot product :  $\oplus$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

## Disentangling Complex Causal Relationships

- Similarly, denote the dynamic soft masks as  $M_t^{\mathcal{D}} \in \mathbb{R}^{|V| \times |V|}$ , to extract dynamic relationships and its complementary to extract the static relationships.
- Dynamic Casual set  $G_{1:T}^{\mathcal{D}} = (X_{1:T}, A_{1:T} \oplus M_{1:T}^{\mathcal{C}} \oplus M_{1:T}^{\mathcal{D}})$
- Static Casual set  $G_{1:T}^{\mathcal{S}} = (X_{1:T}, A_{1:T} \oplus M_{1:T}^{\mathcal{C}} \oplus \bar{M}_{1:T}^{\mathcal{D}})$
- The ground-truth trivial set, dynamic causal set, and static causal set are unavailable in real world applications.
- So, we aim to capture the trivial, dynamic, and static relationships from the full graph by learning the masks.

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Estimating soft mask

- dynamic VGAE-based encoder-decoder to estimate the soft masks of explainable subgraphs.
- At the t-th time step, the causal soft mask matrix can be calculated as

$$\mathbf{M}_t^c = f_v(\mathbf{X}_{1:t}, \mathbf{A}_{1:t}; \Theta_c) = p(\mathbf{M}_t^c | \mathbf{H}_t)q(\mathbf{H}_t | \mathbf{G}_{1:t})$$

- Latent representation calculation

$$q(\mathbf{H}_t | \mathbf{G}_{1:t}) = \prod_{i=1}^N q(\mathbf{h}_{t,i} | \mathbf{G}_{1:t}), q(\mathbf{h}_{t,i} | \mathbf{G}_{1:t}) = \mathcal{N}(\mathbf{h}_{t,i} | \boldsymbol{\mu}_{t,i}, \text{diag}(\boldsymbol{\sigma}_{t,i}^2))$$

- Generating explainable subgraphs:

$$p(\mathbf{M}_t^c | \mathbf{H}_t) = \prod_{i=1}^N \prod_{j=1}^N p(M_{t,ij}^c | \mathbf{h}_{t,i}, \mathbf{h}_{t,j}), p(M_{t,ij}^c = 1 | \mathbf{h}_{t,i}, \mathbf{h}_{t,j}) = g(\mathbf{h}_{t,i}, \mathbf{h}_{t,j})$$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

- static relationship S can also be treated as the co-founder between D and Y, just like the trivial relationship T respect to C and Y.
- The same VGAE-based encoder-decoder framework, different parameters  $\Theta_D$ :

$$\mathbf{M}_t^{\mathcal{D}} = f_v(\mathbf{X}_{1:t}, \mathbf{A}_{1:t} \oplus \mathbf{M}_{1:t}^{\mathcal{C}}; \Theta_D)$$

- Now we have adjacency matrix of Casual Set, Dynamic Causal Set, Static Causal Set
- We need to disentangle the trivial, dynamic, and static relationships



# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Disentangling trivial and causal

- Objective: Ensure Explanation Fidelity
  - Causal Subgraph Set: Target for explanations, representing essential information.
  - Trivial Subgraph Set: Treated as noise, serving as negative examples.
- Approach:
  - Fidelity Criterion: Explanations with causal subgraphs should mimic the original graph's behavior.
  - Negative Treatment: Trivial subgraphs should not influence the model's predictions.

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Disentangling trivial and causal

- Methodology: Embedding Extraction:
- Utilize the pre-trained aggregation function from the DyGNN.
- During each time step  $t$ :
  - $f_a$  would generate the embeddings via extracting the essential information until time  $t$ , for:
  - Original Graph Set  $G_{1:t}$ , Causal Subgraph Set  $G_{1:t}^C$ , Trivial Subgraph Set  $G_{1:t}^T$ .

$$\mathbf{R}_t = f_a(\mathbf{G}_{1:t}), \mathbf{R}_t^C = f_a(\mathbf{G}_{1:t}^C), \mathbf{R}_t^T = f_a(\mathbf{G}_{1:t}^T)$$

- The above outputs can be the node (graph) embedding for downstream tasks.

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Disentangling trivial and causal

- Utilize contrastive learning
- Ensure the semantic similarity between the causal embedding  $e_t^C$  and the original embedding  $e_t$
- Enlarging the semantic distance between the causal embedding  $e_t^C$  and the trivial embedding  $e_t^T$ .

$$\mathcal{L}_c = \frac{1}{T} \sum_{t=1}^T \log \frac{\exp(s(\mathbf{e}_t, \mathbf{e}_t^C)/\tau)}{\exp(s(\mathbf{e}_t, \mathbf{e}_t^C)/\tau) + \alpha_1 \exp(s(\mathbf{e}_t^T, \mathbf{e}_t^C)/\tau) + \alpha_2 \sum_{k \neq t} \exp(s(\mathbf{e}_t^T, \mathbf{e}_k^C)/\tau)}$$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Disentangling static and dynamic

- Extract the dynamic relationship and static relationship
- From the dynamic causal set  $G_{1:T}^D$  and static causal set  $G_{1:T}^S$
- Utilize GCN with learn-able parameters  $\Psi_D$  and  $\Psi_S$ .

$$\mathbf{H}_t^D = GCN(\mathbf{A}_t^D, \mathbf{X}_t; \Psi_D), \mathbf{H}_t^S = GCN(\mathbf{A}_t^S, \mathbf{X}_t; \Psi_S).$$

$$\mathbf{H}_{1:(t-1)}^D \longrightarrow \mathbf{H}_t^D, \quad \mathbf{H}_{1:(t-1)}^S \perp \mathbf{H}_t^S.$$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Disentangling static and dynamic

- we can use the pre-trained aggregation function  $f_a(\cdot)$  again, which extracts the dynamic relationship from the original graph set.
- Fidelity  $\Rightarrow$  Generated dynamic graph set should guarantee that  $f_a$  can extract the dynamic relationship from it
- Dynamic loss:

$$\mathcal{L}_d = \frac{1}{T-1} \sum_{t=2}^T d(f_a(\mathbf{G}_{1:(t-1)}^{\mathcal{D}}), \mathbf{H}_t^{\mathcal{D}})$$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Spatial-temporal explanation

- Due to the highly temporal correlation for dynamic relationships, it would be difficult to disentangle the dynamic relationship.
- Treat the dynamic relationship at time  $t$  as an invention
- Define the causal effect at time  $t$  as follows:

$$\Delta \mathbf{H}_t^{\mathcal{D}} = f_a(\mathbf{G}_{1:t}^{\mathcal{D}}) - f_a(\mathbf{G}_{1:(t-1)}^{\mathcal{D}})$$



# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Spatial-temporal explanation

Combine the causal effect for the dynamic relationship and the static relationship at time  $t$  as the key causal information for  $G_t$

Propose the learn-able weight pooling method to aggregate all the information across all time slots as follows

$$\mathbf{H}_T = \sum_{t=1}^T t_p(\Delta \mathbf{H}_t^{\mathcal{D}} \oplus \mathbf{H}_t^{\mathcal{S}}) \Delta \mathbf{H}_t^{\mathcal{D}} \oplus \mathbf{H}_t^{\mathcal{S}}, \quad t_p(\mathbf{H}) = \text{Softmax}(\Psi_{\mathcal{P}} \mathbf{H} / \|\Psi_{\mathcal{P}}\|)$$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Spatial-temporal explanation

Use aggregated embedding to explain the ground-truth label via prediction loss

$$\mathcal{L}_p = l(f_d(\mathbf{H}_T), \mathcal{Y})$$

Take the sparsity requirement for both the causal graph set and the dynamic causal graph set via the sparsity loss

$$\mathcal{L}_s = \sum_{t=1}^T \frac{\|\mathbf{A}_t^{\mathcal{C}}\|_1 + \|\mathbf{A}_t^{\mathcal{D}}\|_1}{\|\mathbf{A}_t\|_1}$$

# DISENTANGLING COMPLEX CAUSAL RELATIONSHIPS

Summary:

learn the optimal explainable causal subgraphs, dynamic subgraphs, and temporal importance by solving the following optimization problems

$$\min_{\Theta, \Psi} \mathcal{L}(\Theta, \Psi) = \lambda_1 \mathcal{L}_c + \lambda_2 \mathcal{L}_s + \lambda_3 \mathcal{L}_p + \lambda_4 \mathcal{L}_d$$

# EXPERIMENTS



# EXPERIMENTS

- Use 4 synthetic datasets and 2 real-world datasets
- Node classification task and Graph classification task

Dataset	Node classification				Graph classification	
	DBA-Shapes	DTree-Cycles	DTree-Grid	Elliptic	DBA-2motifs	MemeTracker
#nodes	700	871	1,231	203,769	25,000	3.3 mil.
#edges	4,110	1,950	3410	234,355	51,392	27.6 mil.
#labels	7	3	3	2	3	2

Table 2: Explanation accuracy of different models (%). Where best performances are bold.

Task	Dataset	GNNExplainer	PGExplainer	Gem	OrphicX	DyGNNExplainer
Node cls.	DBA-Shapes	92.1	92.9	93.6	94.3	<b>97.8*</b>
	DTree-Cycles	92.8	93.7	94.4	96	<b>98.2*</b>
	DTree-Grid	85.2	85.9	87.1	90.5	<b>94.2*</b>
	Elliptic	92.4	94.1	94.6	96.1	<b>98.7*</b>
Graph cls.	DBA-2motifs	86.5	88.0	90.7	91.4	<b>96.3*</b>
	MemeTracker	88.2	89.2	91.0	91.9	<b>97.4*</b>

“\*” indicates the statistically significant improvements (i.e., two-sided t-test with  $p < 0.05$ ) over the best baseline. ‘cls.’ is short for classification.

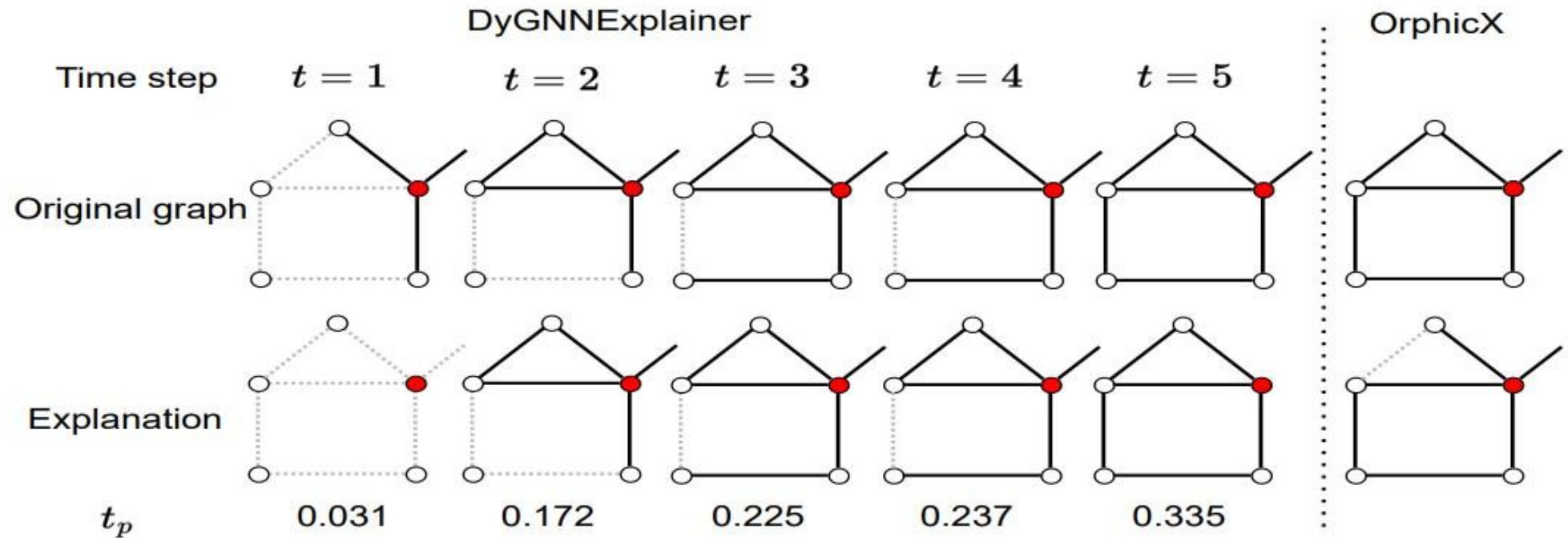
# EXPERIMENTS

- Explanation fidelity: We compare the predicted labels of explanatory subgraphs with the predicted labels of the original graphs as generated by the target model.
- Explanation interpretability analysis: The explanation subgraphs should exhibit a high degree of sparsity. Measure the number of subgraph edges

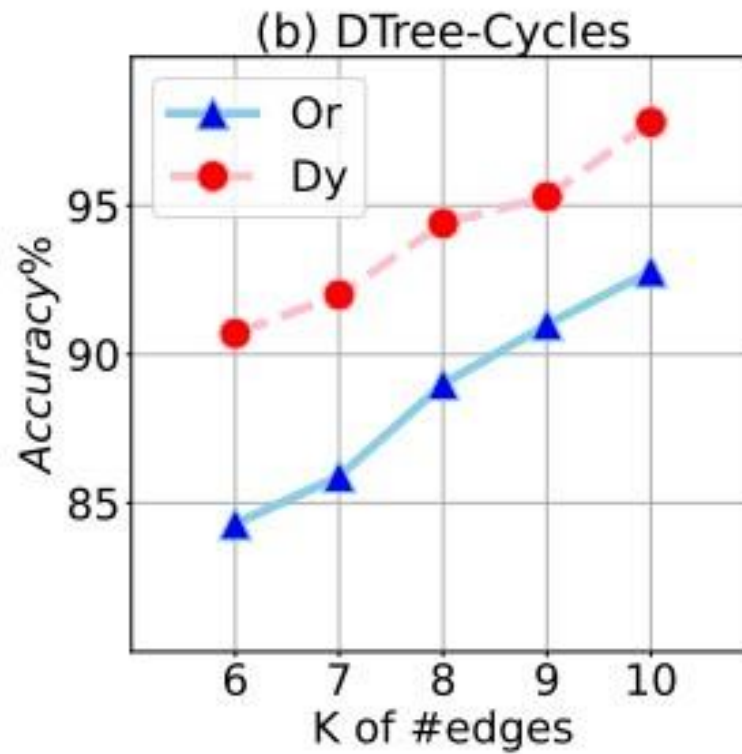
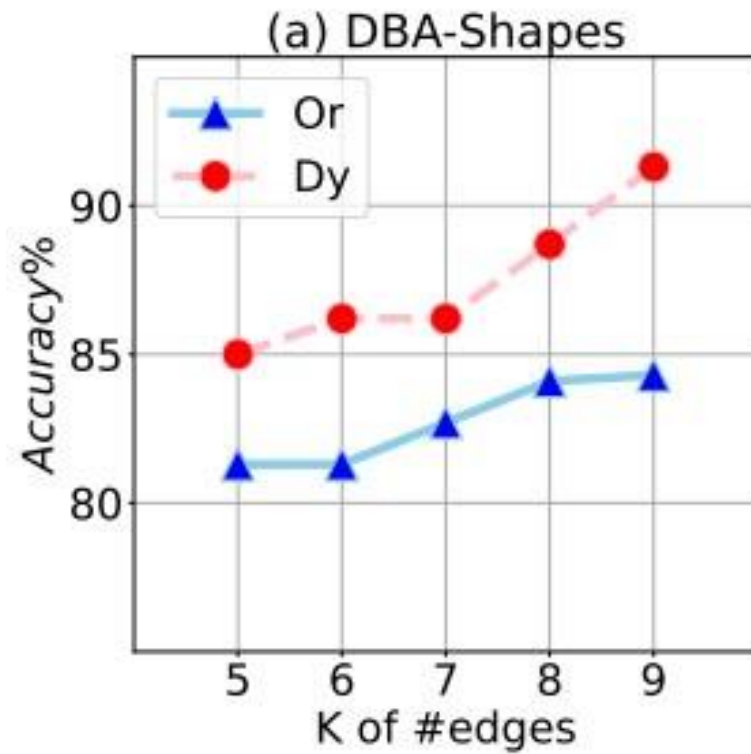


# EXPERIMENTS

## Case Study



# EXPERIMENTS



# EXPERIMENTS

- Prediction accuracy analysis

Table 3: Prediction accuracy of different models (%). Where best performances are bold.

Dataset	GNNExplainer	PGExplainer	Gem	OrphicX	Target	DyGNNExplainer
DBA-Shapes	35.5	36.3	38.5	38.7	40.2	<b>44.6*</b>
Elliptic	39.7	45.6	43.5	47.8	84.3	<b>89.2*</b>

\*,” indicates the statistically significant improvements (i.e., two-sided t-test with  $p < 0.05$ ) over the best baseline.

# RELATED WORK

- Recent methods have emerged to provide explanations for GNNs
- Difference between this work and them

# RELATED WORK

- Methods predominantly aim to generate input-dependent explanations
- GNNExplainer (Ying et al., 2019) seeks soft masks for edges and node features through mask optimization to explain predictions
- Typically explain each instance individually and lack the ability to generalize graphs
- PGExplainer (Luo et al., 2020) proposes learning a mask predictor for edge masks to provide explanations.
- XGNN (Yuan et al., 2020) focuses on investigating graph patterns leading to specific classes
- In contrast to these approaches, this work leverages **causality** to achieve faithful explanations

- Addressed the critical challenges associated with interpretability in DyGNNs.
- Superior performance of DyGNNExplainer in both explanation tasks and real predictions
- Generated synthetic dynamic datasets tailored for dynamic graph interpretability tasks for similar future work

# CONCLUSION





# THANK YOU

Hooman Zolfaghari