

The background of the slide is a light gray gradient with several realistic water droplets of various sizes scattered across it. The droplets have highlights and shadows, giving them a three-dimensional appearance. The title text is centered in a purple color.

# Deep Learning with Graph-Structured Data: A Mathematical Perspective

Jie Chen

MIT-IBM Watson AI Lab, IBM Research

Presented at 47th Annual Spring Lecture Series, University of Arkansas

# About me /



The [MIT-IBM Watson AI Lab](#) is an academic-industrial collaboration model dedicated to pushing the frontiers of artificial intelligence and translating fundamental scientific breakthroughs into business and real-world impacts.

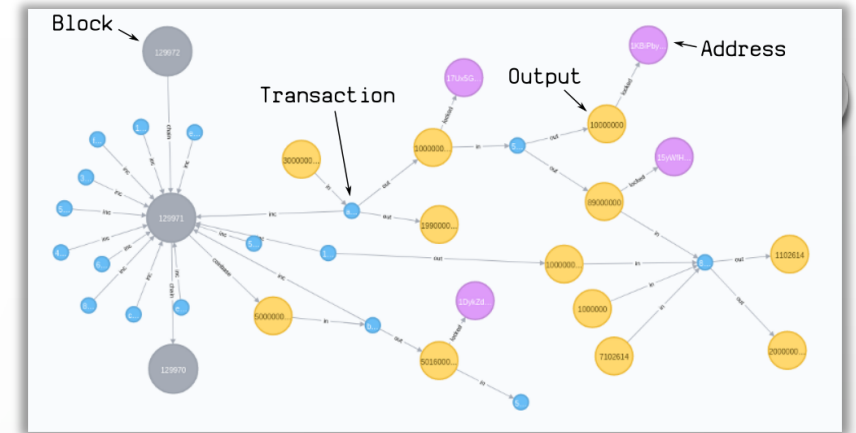




Power grid



Social network



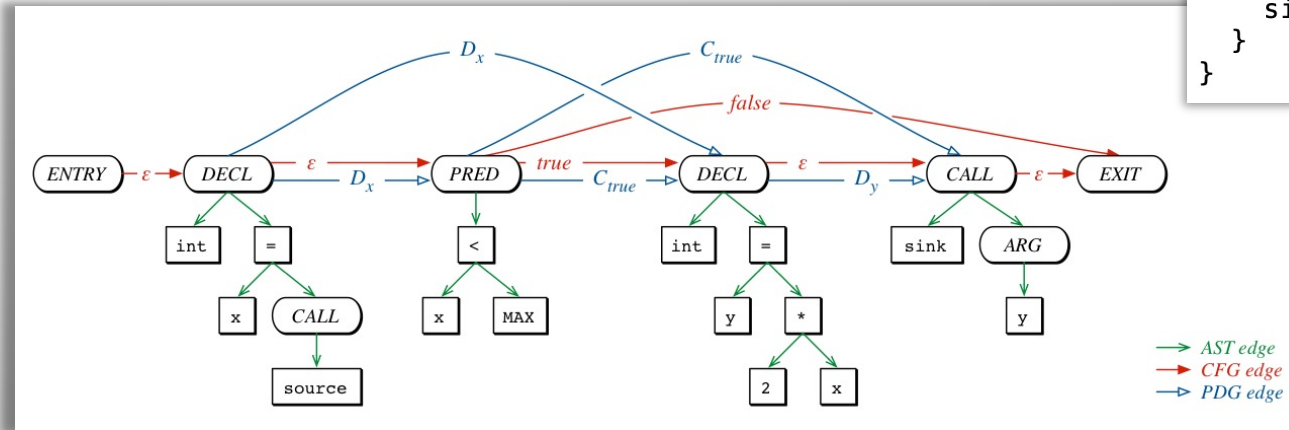
Bitcoin transactions

# Graphs are ubiquitous

Molecules



Code graphs

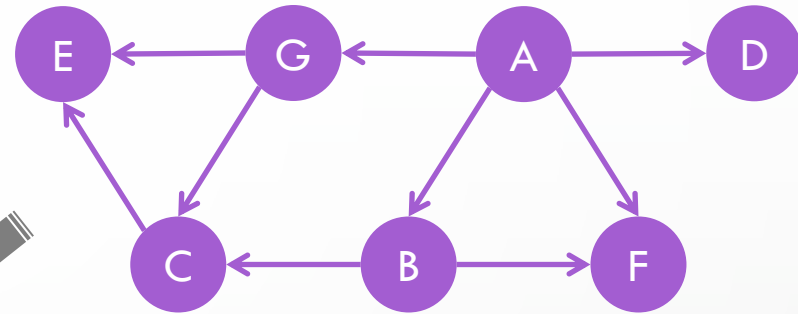


```

void foo() {
  int x = source();
  if (x < MAX) {
    int y = 2 * x;
    sink(y);
  }
}
  
```

→ AST edge  
 → CFG edge  
 → PDG edge

# Sparse matrices $\Leftrightarrow$ graphs



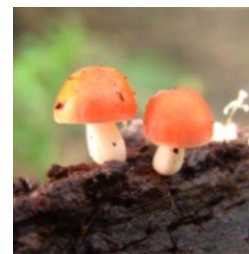
adjacency matrix

	A	B	C	D	E	F	G
A		1		1		1	1
B			1			1	
C					1		
D							
E							
F							
G			1		1		

reordering  
(topological sort)

	A	B	D	F	G	C	E
A		1	1	1	1		
B				1		1	
D							
F							
G						1	1
C							1
E							

# Machine learning 101



$f \rightarrow$  agaric (0.88)  
jelly fungus (0.11)  
fingers (0.01)

- **Model**  $f(x; \theta)$ 
  - return a probability vector (length  $c$ )
  - model parameters
  - input data
  - softmax( $x$ ) =  $(e^{x^1}, e^{x^2}, \dots, e^{x^c})$ /normalization

- **Example**  $f(x; \theta) = \text{softmax}(W^1 \cdot \text{ReLU}(W^0 x + b^0) + b^1)$ ,  $\theta = \{W^0, W^1, b^0, b^1\}$

$$\text{ReLU}(x) = \max(x, 0)$$

$y_i$  is the labeling vector (one-hot)

- **Loss**  $L(\theta) = \frac{1}{n} \sum_{i=1}^n L_i(\theta)$ ,  $L_i(\theta) = -(y_i)_1 \log f(x_i; \theta)_1 - \dots - (y_i)_c \log f(x_i; \theta)_c$

de-facto choice: stochastic optimization

- **Training**  $\min_{\theta} L(\theta)$

$$\text{Mini-batch loss } L_B(\theta) = \frac{1}{|B|} \sum_{i \in B} L_i(\theta) ; \quad \text{Update rule: } \theta^{k+1} = \theta^k - \gamma^k \nabla L_B(\theta^k)$$

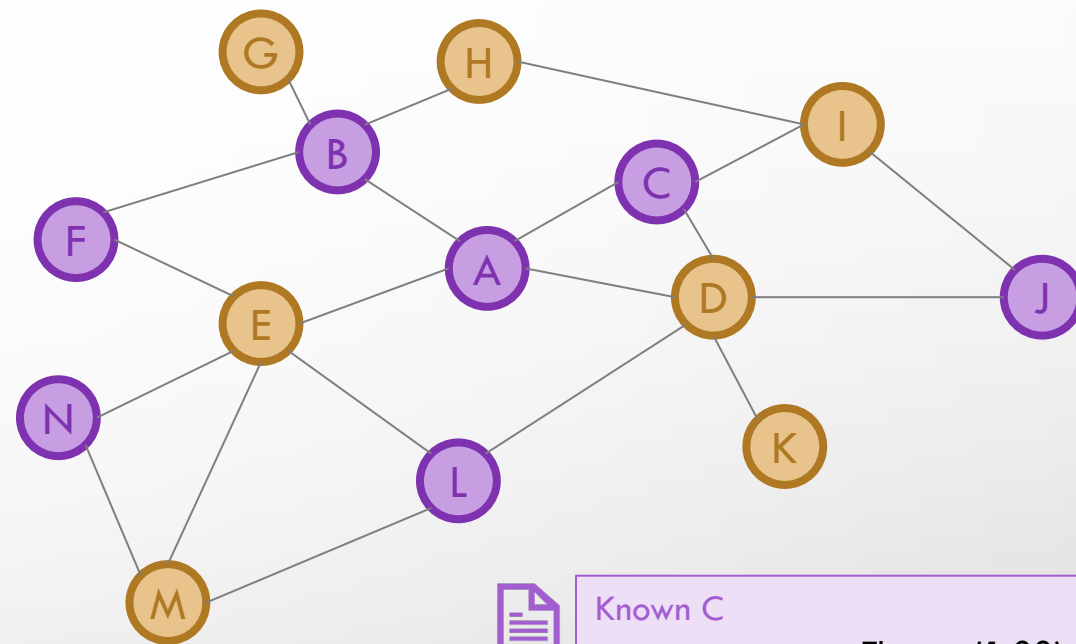
# Graph data and task

From machine learning 101:

- A data point is  $x_i \in \mathbb{R}^d$  with label  $y_i \in \mathbb{R}^c$
- Some labels are known, while some not

Additionally:

- Data points are connected by a (weighted) adjacency matrix  $A$
- The task is to predict the unknown labels



Known C

Theory (1.00)  
Neural networks (0.00)  
Reinforcement learning (0.00)

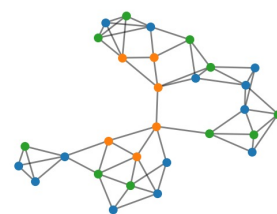


Predict D

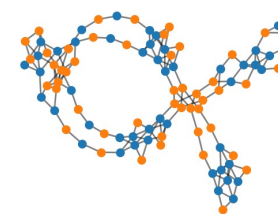
Theory (?)  
Neural networks (?)  
Reinforcement learning (?)

# More graph data and tasks

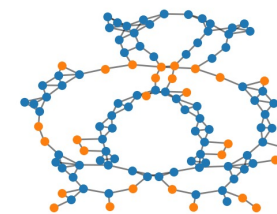
- **Graph-level prediction:** The entire graph is a data point  $x_i$ ; the task is to predict the label  $y_i$  of this graph
- **Time series forecasting:** Split time interval  $T = T_1 + T_2$ , where input  $x_i \in \mathbb{R}^{T_1}$  and output  $y_i \in \mathbb{R}^{T_2}$ . Additionally, several time series are interconnected by a graph.



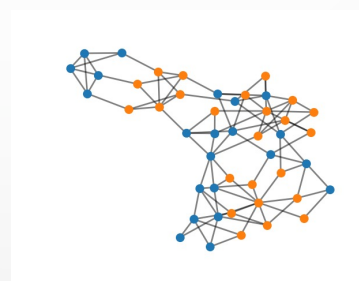
enzyme



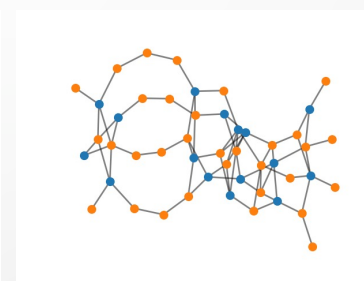
enzyme



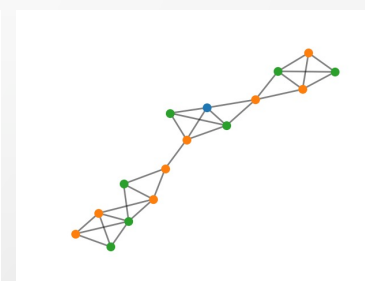
non-enzyme



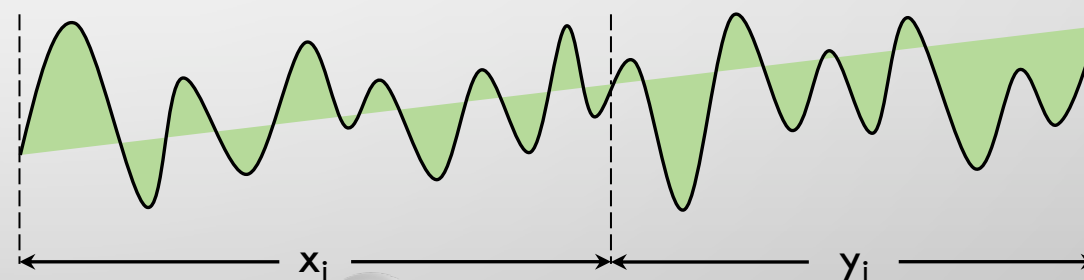
?



?



?



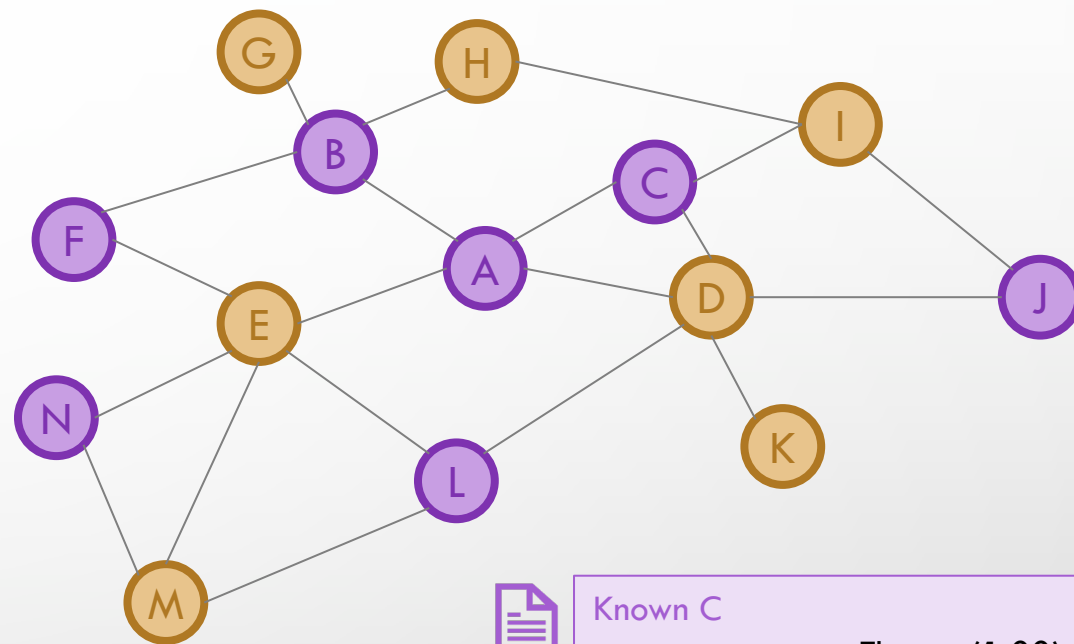
# Graph neural networks (GNN)

- Graph convolution network (GCN)

$$f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$$

$$f(X) = \begin{matrix} \text{purple} \\ \text{orange} \\ \text{purple} \end{matrix} \quad X = \begin{matrix} \text{purple} \\ \text{orange} \\ \text{purple} \end{matrix}$$

$$\hat{A} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}, \quad \tilde{A} = A + I, \quad \tilde{D} = \text{diag} \left( \sum_j \tilde{A}_{ij} \right)$$



Known C

Theory (1.00)

Neural networks (0.00)

Reinforcement learning (0.00)



Predict D

Theory (?)

Neural networks (?)

Reinforcement learning (?)



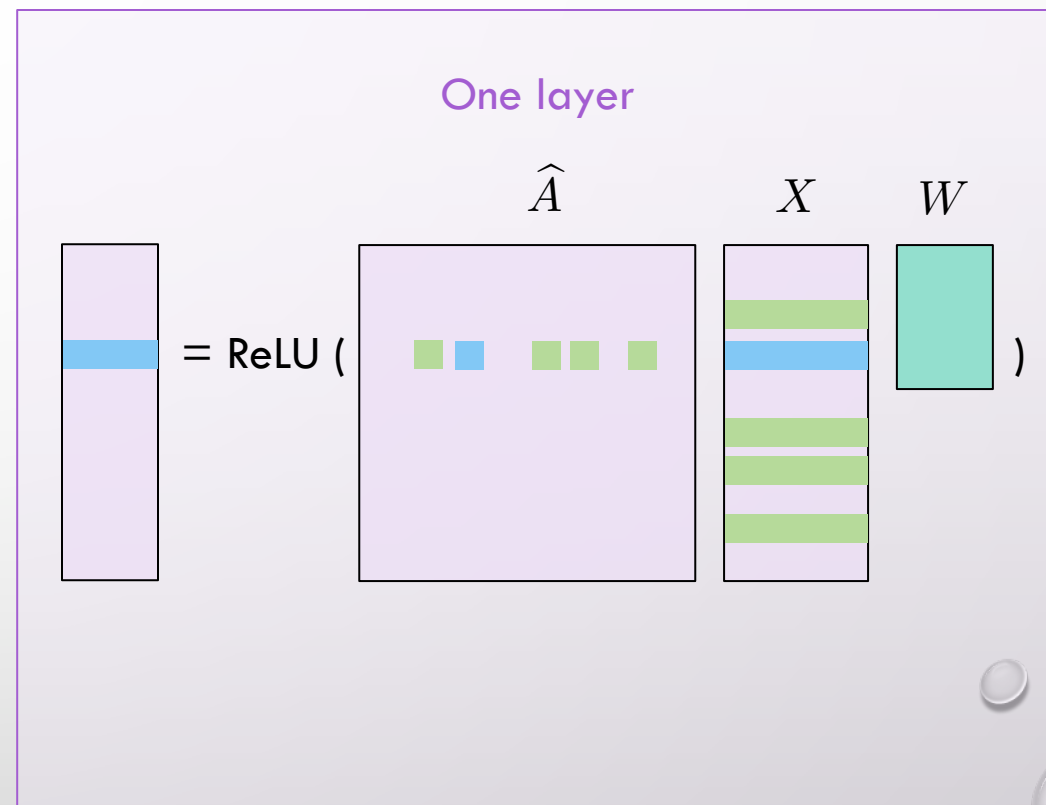
# Graph neural networks (GNN) --- Matrix view

- Graph convolution network (GCN)

$$f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$$

$$f(X) = \begin{matrix} \text{purple} \\ \text{orange} \\ \text{purple} \\ \text{orange} \end{matrix} \quad X = \begin{matrix} \text{purple} \\ \text{orange} \\ \text{purple} \\ \text{orange} \end{matrix}$$

$$\hat{A} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}, \quad \tilde{A} = A + I, \quad \tilde{D} = \text{diag} \left( \sum_j \tilde{A}_{ij} \right)$$



# Graph neural networks (GNN) --- Message passing

- Graph convolution network (GCN)

$$f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$$

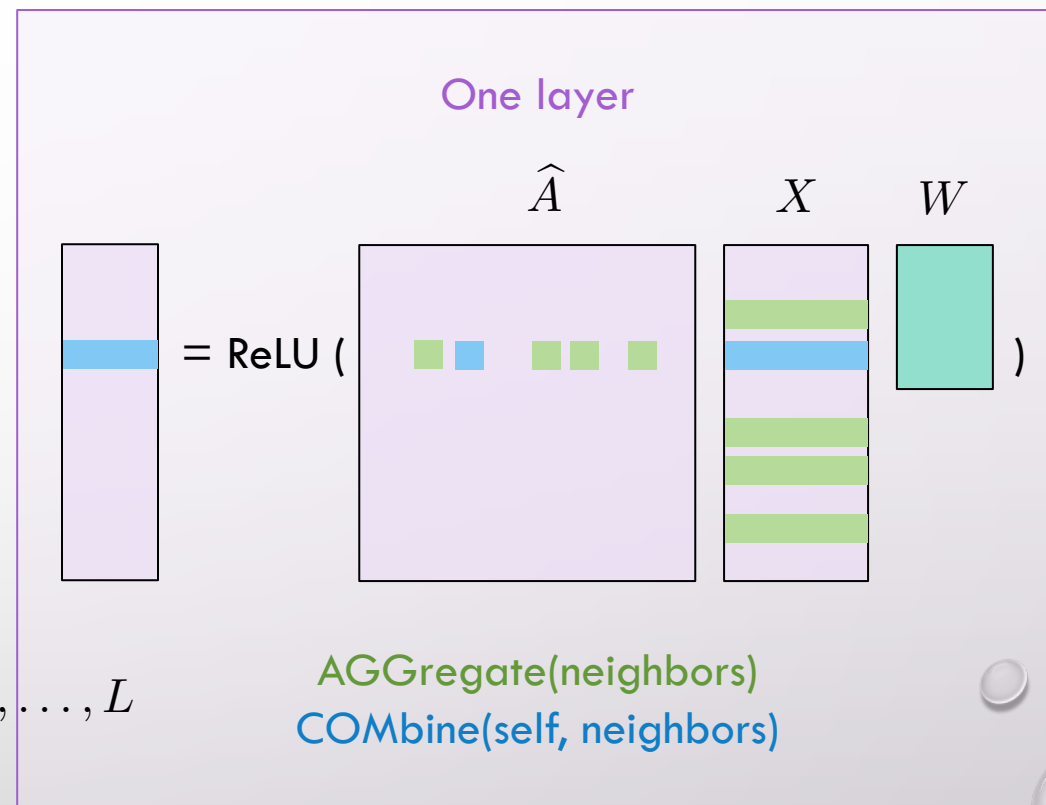
$$f(X) = \begin{bmatrix} \text{purple} \\ \text{orange} \\ \text{purple} \\ \text{orange} \end{bmatrix} \quad X = \begin{bmatrix} \text{purple} \\ \text{orange} \\ \text{purple} \\ \text{orange} \end{bmatrix}$$

$$\hat{A} = \tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}}, \quad \tilde{A} = A + I, \quad \tilde{D} = \text{diag} \left( \sum_j \tilde{A}_{ij} \right)$$

- Message-passing framework

$$h_i^\ell = \text{COM}^\ell \left( h_i^{\ell-1}, \text{AGG}^\ell(\{h_j^{\ell-1} \mid j \in \mathcal{N}_i\}) \right), \quad \ell = 1, \dots, L$$

$$f(x_i) = h_i^L \quad x_i = h_i^0$$



## More examples of message passing

For simplicity, drop superscript (layer index):  $h'_i = \text{COM}(h_i, \text{AGG}(\{h_j \mid j \in \mathcal{N}_i\}))$

### GraphSAGE

$$h'_i = \text{ReLU} \left( W_1 h_i + W_2 \cdot \text{mean}_{j \in \mathcal{N}_i} h_j \right)$$

### Graph isomorphism network (GIN)

$$h'_i = \text{NeuralNetwork} \left( (1 + \epsilon) \cdot h_i + \sum_{j \in \mathcal{N}_i} h_j \right)$$

### Graph attention network (GAT)

$$h'_i = \text{ReLU} \left( \alpha_{ii} W h_i + \sum_{j \in \mathcal{N}_i} \alpha_{ij} W h_j \right)$$
$$\alpha_{ij} = \frac{\exp(\text{LeakyReLU}(a_1^T W h_i + a_2^T W h_j))}{\text{normalize over } \{j \in \{i\} \cup \mathcal{N}_i\}}$$

## Questions to consider

All models are wrong;  
some are useful.  
– George Box

1. Why do GNNs work?
2. Do GNNs scale to massive graphs?
3. What if the graph evolves over time?
4. Are GNNs robust to perturbations of the graph?
5. Are there better-performing GNNs for special types of graphs (e.g., DAG)?
6. What if the weight matrices  $W$  are infinitely large? What if the number of layers,  $L$ , is infinite?
7. If there is not a graph, can we learn one so that we can apply GNNs?
8. How can we generate new graphs?



# Learning with massive graphs

**References:** Chen, Ma, and Xiao (2018). FastGCN: Fast Learning with Graph Convolutional Networks via Importance Sampling.

Chen and Luss (2018). Stochastic Gradient Descent with Biased but Consistent Gradient Estimators.



## Recall the setting

- $n$  nodes,  $d$  input features,  $c$  classes
- Train a model  $f$  such that  $f(X) \approx Y$ ,  $X \in \mathbb{R}^{n \times d}$ ,  $Y \in \mathbb{R}^{n \times c}$ ,  $A \in \mathbb{R}^{n \times n}$
- Focus on the GCN model  $f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$
- Question: How do we find the parameters of  $f$  (e.g.,  $W^0$  and  $W^1$ ) efficiently for large  $n$ ?

# Challenge

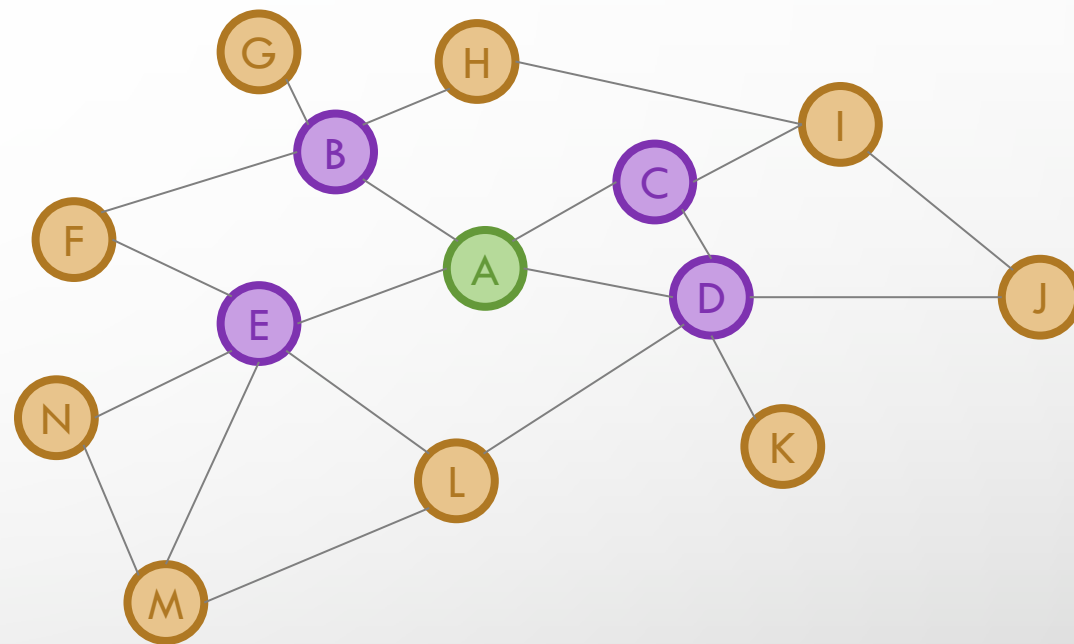
$$f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$$

one or a  
small batch

need  
neighbors

need neighbors'  
neighbors

- In stochastic optimization, each step randomly picks one (or a small batch of) node(s) to compute the approximate gradient
- For GCN, the worst-case cost is  $O(n)$  per gradient step, for small-world graphs



# Solution: Sampling

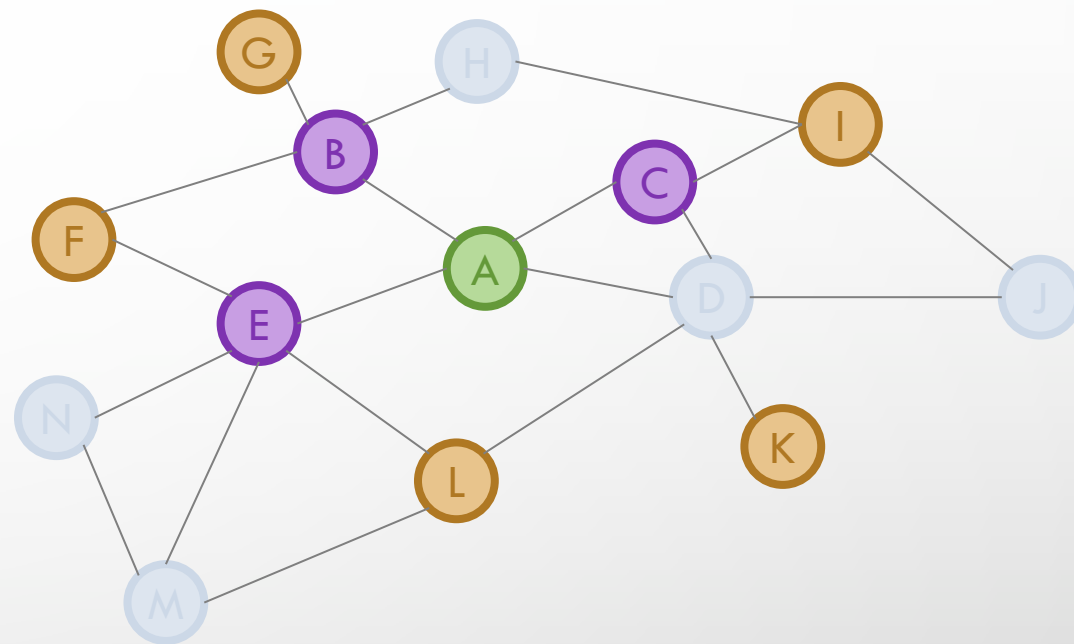
$$f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$$

one or a  
small batch

3 one-hop  
neighbors

5 two-hop  
neighbors

- Randomly pick (according to any probability distribution) a fixed number of neighbors only, for each layer
- The worst-case cost is  $O(1)$  per gradient step





# Theory: Sample gradient is not unbiased, but consistent

- Assume objective  $L(\theta)$  is nonconvex but  $L$ -smooth; gradient is bounded with  $\|\nabla L(\theta)\| \leq G$
- Use stochastic gradient update  $\theta_{k+1} = \theta_k - \gamma_k g_k$  where  $g_k$  is consistent
- Specifically,  $\Pr\left(\|g_k - \nabla L(\theta_k)\| \geq \delta \|\nabla L(\theta_k)\| \mid g_1, \dots, g_{k-1}\right) \leq C_k e^{-N_k \tau(\delta)}$  where  $\delta$  is in  $(0,1)$  and  $\tau$  is an increasing and positive function
- Use constant (oracle) step size  $\gamma_k = D_f / [(1 + \delta)G\sqrt{T}]$  where  $D_f = [2(L(\theta_1) - L(\theta^*)) / L]^{\frac{1}{2}}$
- For any  $\epsilon$  in  $(0,1)$ , using sample size  $N_k \geq \tau(\delta)^{-1} \log(TC_k/\epsilon)$ , with probability  $\geq 1 - \epsilon$ ,

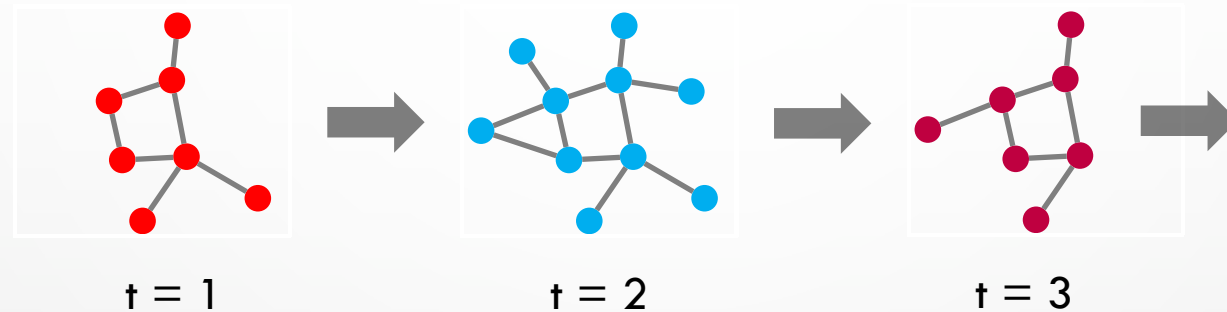
$$\min_{k=1, \dots, T} \|\nabla L(\theta_k)\|^2 \leq \frac{(1 + \delta) L G D_f}{(1 - \delta) \sqrt{T}}$$



# Learning with dynamic graphs

**References:** Pareja, et al (2020). EvolveGCN: Evolving Graph Convolutional Networks for Dynamic Graphs.

## In most real-life applications, graphs are dynamic



- $n$  nodes,  $d$  input features,  $c$  classes,  $T$  time steps
- Train a model  $f$  such that  $f(X_t) \approx Y_t$ ,  $X_t \in \mathbb{R}^{n \times d}$ ,  $Y_t \in \mathbb{R}^{n \times c}$ ,  $A_t \in \mathbb{R}^{n \times n}$ ,  $t = 1, \dots, T$
- The number of nodes,  $n$ , can change over time (but for simplicity we call it  $n$ )
- Question: How do we build such a model based on one for static graphs?

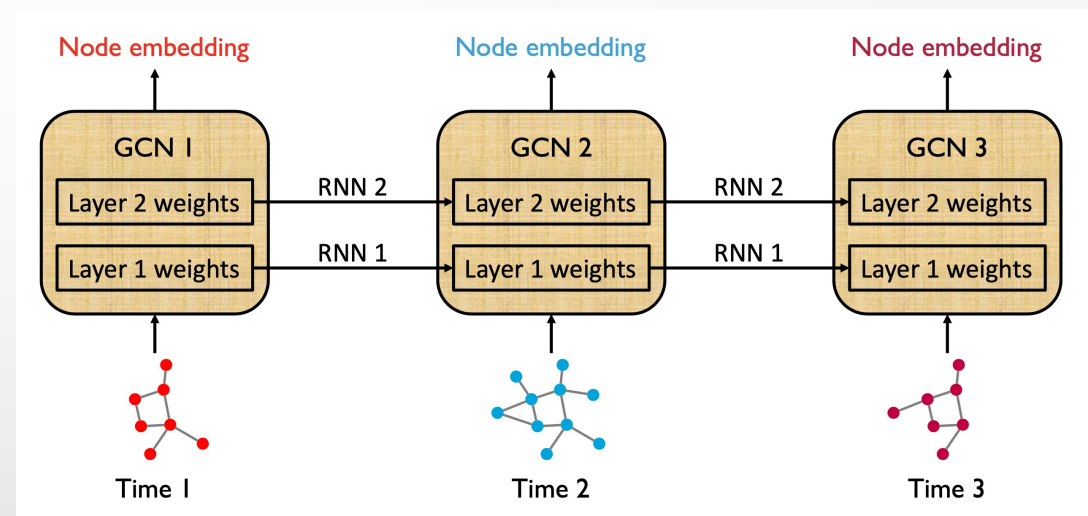
## A closer look at GCN

$$f(X) = \text{softmax}(\hat{A} \cdot \text{ReLU}(\hat{A}XW^0) \cdot W^1)$$

layer notation:  $H^{\ell+1} = \sigma(\hat{A}H^{\ell}W^{\ell})$

When  $A$  evolves over time, we have two choices

- Model the evolution of layer embeddings  $H^{\ell}$  by using a recurrent neural network (RNN)
  - This is a typical approach, but it cannot handle the dynamic appearance/disappearance of nodes
- Model the evolution of layer parameters  $W^{\ell}$  by using a recurrent neural network (RNN)



## Two evolution mechanisms

- H version: GCN parameters are GRU hidden states

$$W_t^\ell = \text{GRU}(H_t^\ell, W_{t-1}^\ell)$$

$$H_t^{\ell+1} = \sigma(\hat{A}_t H_t^\ell W_t^\ell)$$

- O version: GCN parameters are LSTM input/outputs

$$W_t^\ell = \text{LSTM}(W_{t-1}^\ell)$$

$$H_t^{\ell+1} = \sigma(\hat{A}_t H_t^\ell W_t^\ell)$$

1: **function**  $H_t = \text{GRU}(X_t, H_{t-1})$

// Input  $X_t$ , hidden states  $H_{t-1}$  and  $H_t$

2:  $Z_t = \text{sigmoid}(W_Z X_t + U_Z H_{t-1} + B_Z)$

3:  $R_t = \text{sigmoid}(W_R X_t + U_R H_{t-1} + B_R)$

4:  $\tilde{H}_t = \tanh(W_H X_t + U_H (R_t \circ H_{t-1}) + B_H)$

5:  $H_t = (1 - Z_t) \circ H_{t-1} + Z_t \circ \tilde{H}_t$

6: **end function**

1: **function**  $H_t = \text{LSTM}(X_t)$

// Current input  $X_t =$  past output  $H_{t-1}$

2:  $F_t = \text{sigmoid}(W_F X_t + U_F H_{t-1} + B_F)$

3:  $I_t = \text{sigmoid}(W_I X_t + U_I H_{t-1} + B_I)$

4:  $O_t = \text{sigmoid}(W_O X_t + U_O H_{t-1} + B_O)$

5:  $\tilde{C}_t = \tanh(W_C X_t + U_C H_{t-1} + B_C)$

6:  $C_t = F_t \circ C_{t-1} + I_t \circ \tilde{C}_t$

7:  $H_t = O_t \circ \tanh(C_t)$

8: **end function**

# Parameter evolution compared with other static/dynamic models

Table 2: Performance of link prediction. Each column is one data set.

	mean average precision					mean reciprocal rank				
	SBM	BC-OTC	BC-Alpha	UCI	AS	SBM	BC-OTC	BC-Alpha	UCI	AS
GCN	0.1987	0.0003	0.0003	0.0251	0.0003	0.0138	0.0025	0.0031	0.1141	0.0555
GCN-GRU	0.1898	0.0001	0.0001	0.0114	0.0713	0.0119	0.0003	0.0004	0.0985	0.3388
DynGEM	0.1680	0.0134	0.0525	0.0209	0.0529	0.0139	0.0921	0.1287	0.1055	0.1028
dyngraph2vecAE	0.0983	0.0090	0.0507	0.0044	0.0331	0.0079	0.0916	0.1478	0.0540	0.0698
dyngraph2vecAERNN	0.1593	<b>0.0220</b>	<b>0.1100</b>	0.0205	0.0711	0.0120	<b>0.1268</b>	<b>0.1945</b>	0.0713	0.0493
EvolveGCN-H	0.1947	0.0026	0.0049	0.0126	<b>0.1534</b>	<b>0.0141</b>	0.0690	0.1104	0.0899	<b>0.3632</b>
EvolveGCN-O	<b>0.1989</b>	0.0028	0.0036	<b>0.0270</b>	0.1139	0.0138	0.0968	0.1185	<b>0.1379</b>	0.2746

Our method (parameter evolution)



# Learning to construct a graph

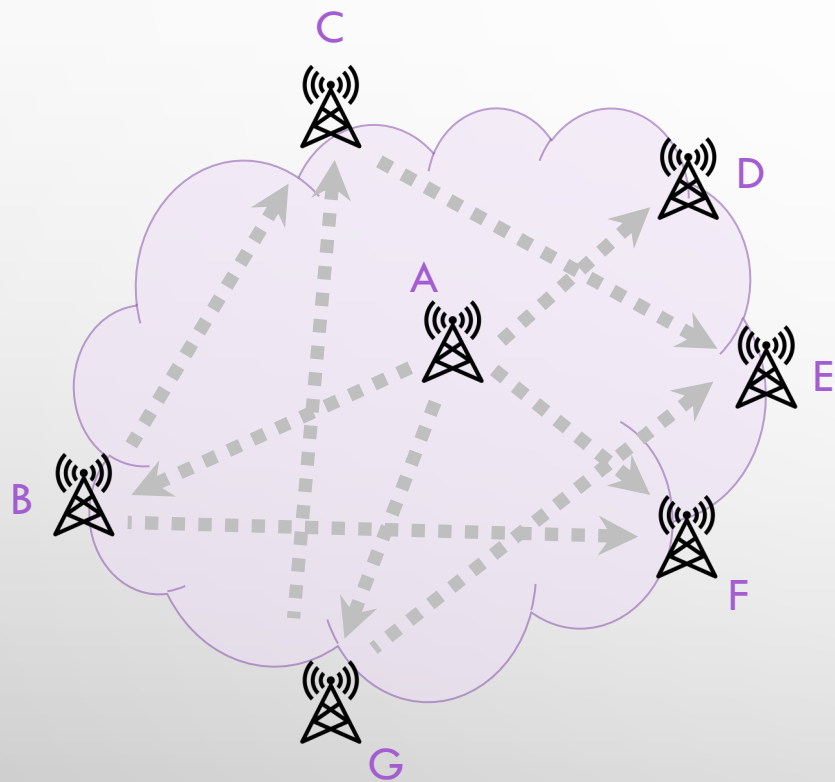
**References:** Shang, Chen, and Bi (2021). Discrete Graph Structure Learning for Forecasting Multiple Time Series.

Dai and Chen (2022). Graph-Augmented Normalizing Flows for Anomaly Detection of Multiple Time Series.

Yu, et al (2019). DAG-GNN: DAG Structure Learning with Graph Neural Networks.



# Multiple time series



A



B



C



D



E



**Forecasting task:** Predicting values for the next few time steps

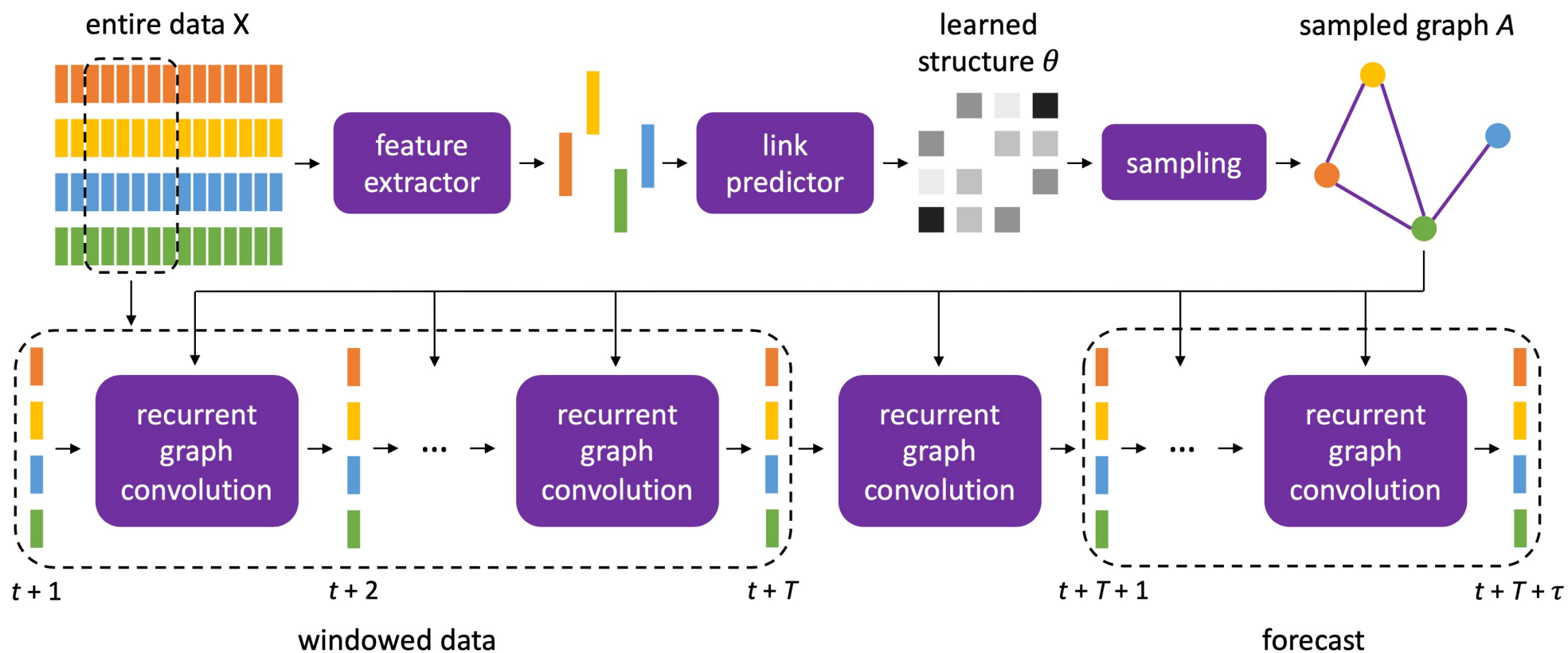
**Anomaly detection task:** Find out time series that behave abnormally

Sensors are connected, leading to interactions/interdependency

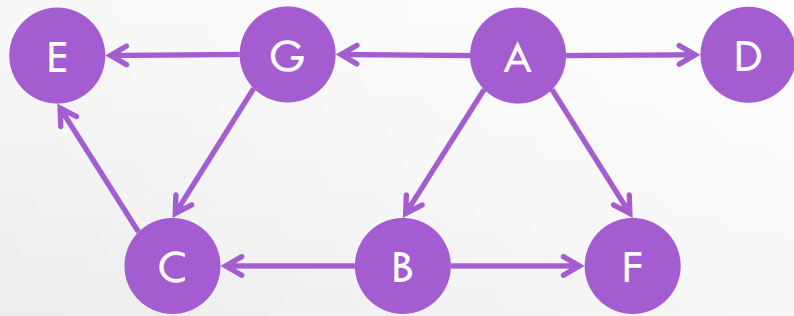
- Can we learn the interactions/interdependency?
- Do they help the task?



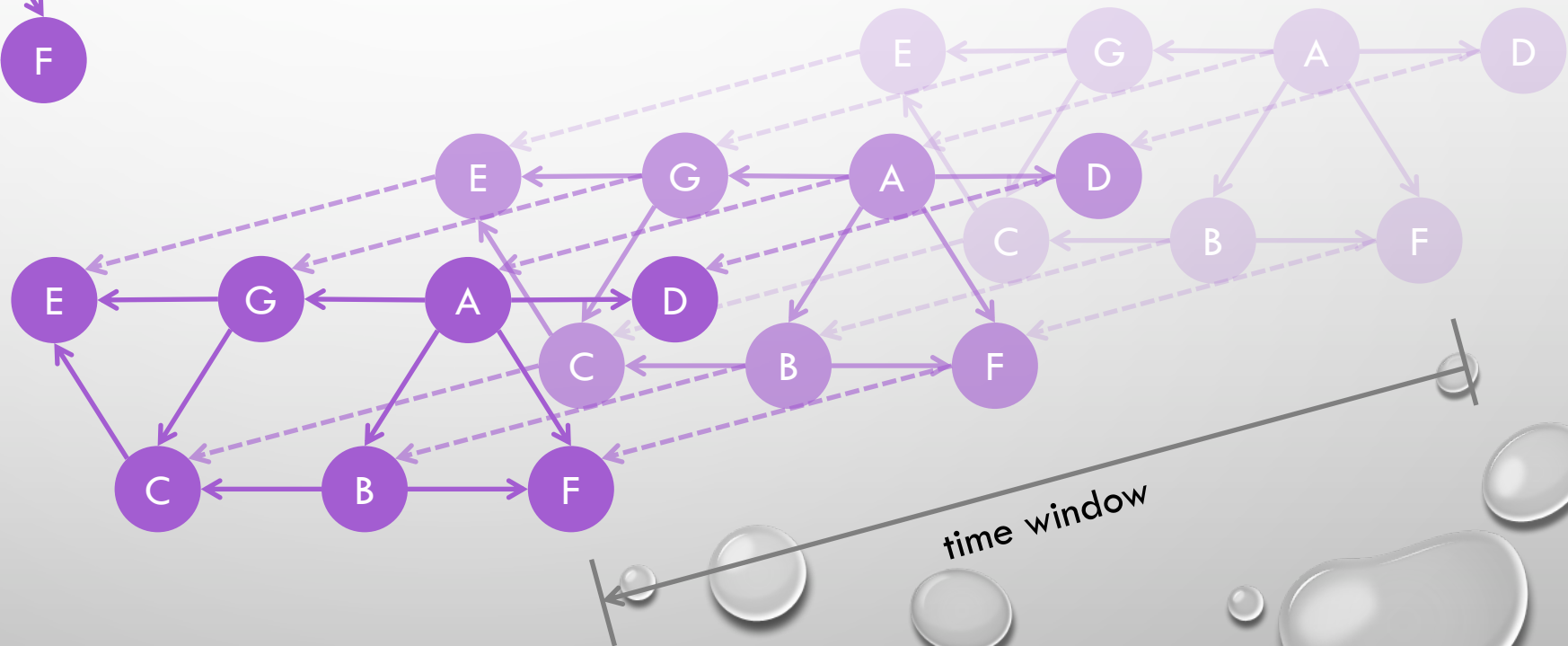
# A link prediction approach



# A Bayesian network approach (directed acyclic)



$$P(A, B, C, D, E, F, G) = P(E|C, G) \cdot P(G|A) \cdot P(A) \cdot P(D|A) \cdot P(C|B, G) \cdot P(B|A) \cdot P(F|A, B)$$



# A Bayesian network approach

$$P(A_1, B_1, C_1, D_1, E_1, F_1, G_1, A_2, B_2, C_2, D_2, E_2, F_2, G_2, A_3, B_3, C_3, D_3, E_3, F_3, G_3)$$

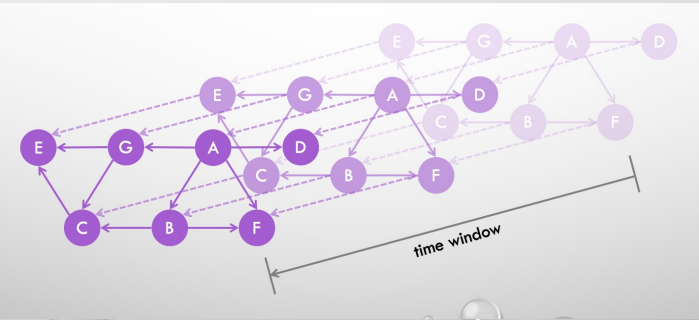
$$= P(A_1) * P(A_2 | A_1) * P(A_3 | A_2) * P(B_1 | A_1) * P(B_2 | A_2, B_1) * P(B_3 | A_3, B_2) * P(C_1 | B_1, G_1) * P(C_2 | B_2, G_2, C_1) * P(C_3 | B_3, G_3, C_2) * P(D_1 | A_1) * P(D_2 | A_2, D_1) * P(D_3 | A_3, D_2) * P(E_1 | C_1, G_1) * P(E_2 | C_2, G_2, E_1) * P(E_3 | C_3, G_3, E_2) * P(F_1 | A_1, B_1) * P(F_2 | A_2, B_2, F_1) * P(F_3 | A_3, B_3, F_2) * P(G_1 | A_1) * P(G_2 | A_2, G_1) * P(G_3 | A_3, G_2)$$

normalizing flow



$$P(B_3 | A_3, B_2)$$

dependency encoder; graph convolution



# Yes, learning a graph indeed helps

Table 2: Forecasting error (PMU).

	Metric	FNN	LSTM	DCRNN	LDS	NRI	GTS <sub>v</sub>	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>

Our method

Table 1: AUC-ROC (%) of anomaly detection.

Dataset	EncDecAD	DeepSVDD	ALOCC	DROCC	DeepSAD	GANF
PMU-B	55.6 $\pm$ 1.8	55.6 $\pm$ 3.3	62.9 $\pm$ 2.2	58.6 $\pm$ 3.0	63.7 $\pm$ 0.9	<b>67.5<math>\pm</math>0.8</b>
PMU-C	53.7 $\pm$ 0.5	56.9 $\pm$ 0.9	60.9 $\pm$ 1.3	61.9 $\pm$ 2.7	60.1 $\pm$ 1.4	<b>70.6<math>\pm</math>3.3</b>
SWaT	76.5 $\pm$ 0.7	68.8 $\pm$ 2.0	75.4 $\pm$ 2.3	73.3 $\pm$ 1.6	75.4 $\pm$ 1.2	<b>79.6<math>\pm</math>0.9</b>

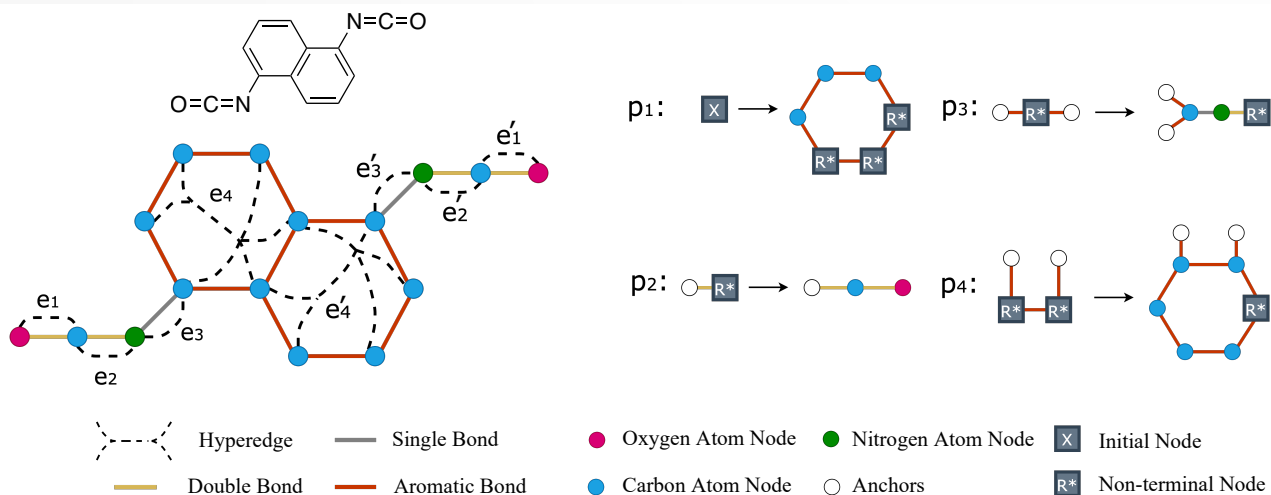
Our method

# Epilogue

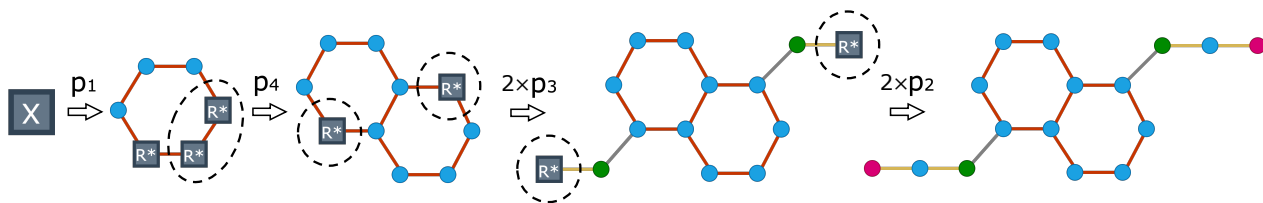


# Some graph projects at MIT-IBM

## Learning a graph grammar for molecule generation and optimization

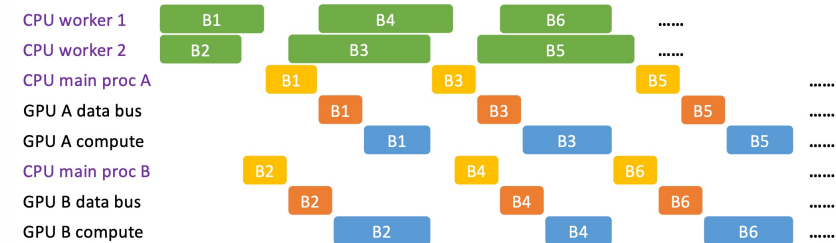


(a) Hypergraph representation of naphthalene diisocyanate (b) Production rules of a learned graph grammar

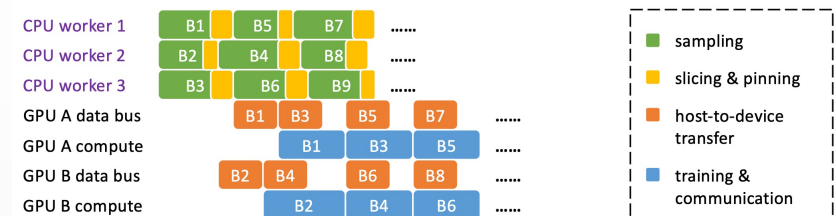


(c) Generation process of naphthalene diisocyanate using the graph grammar in (b)

## Training and inference system (in-memory/out-of-core)



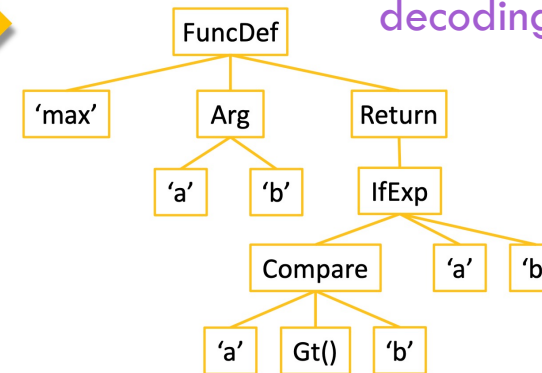
(a) Standard PyTorch workflow.



(b) Our system SALIENT.

```
def max(a, b):
    return a if a > b else b
```

Program translation  
through encoding-  
decoding parse trees



```
public static int max(int a, int b) {
    return (a > b) ? a : b;
}
```