

# Graph Representation Learning

**Tianxiang (Adam) Gao**

November 11, 2024

# Outline

- 1 Machine Learning with Graphs
- 2 Introduction to Graphs
- 3 Graph Neural Networks
- 4 Training GNNs

# Recap: Transformers

- **Self-Attention:** Refines the representation of each token by learning its relevance to all other tokens, *i.e.*,  $z = \sum_i \alpha_i v_i$ , where  $\alpha_i$  represents attention weights.
- **Multi-Head Attention:** Focuses on different aspects of each token to capture diverse patterns, *i.e.*,  $z = [z_1 \dots z_H] W_o$ , where each  $z_h$  represents an individual attention head.
- **Layer Normalization:** Normalizes each layer by computing statistics across the hidden units within a layer.
- **Encoder-Decoder Attention:** Refines the output representation by referencing the input representations.
- **Masked Attention:** Masks future tokens to maintain autoregressive generation, preventing “leakage” of future information.
- **Positional Encoding:** Provides unique, low-dimensional representations to encode token positions, allowing the model to differentiate positional relationships easily.
- **Teacher Forcing:** Uses the correct prior output during the training to facilitate learning.

## Recap: Large Language Models

- **BERT**: An encoder-only architecture that utilizes both past and future context in bidirectional self-attention.
- **Masked Language Modeling**: BERT is pretrained by predicting masked tokens based on surrounding context.
- **GPT**: A decoder-only architecture, pretrained as a standard language model that predicts the next token in a sequence.
- **Zero-Shot Inference**: The pretrained model performs inference by prompting with a task description without fine-tuning.
- **In-Context Learning**: The model learns to perform tasks by providing demonstrations before posing the question.
- **Neural Scaling Laws**: Predicts computational efficiency gains as model size increases, showing systematic improvement with model scaling.

# Outline

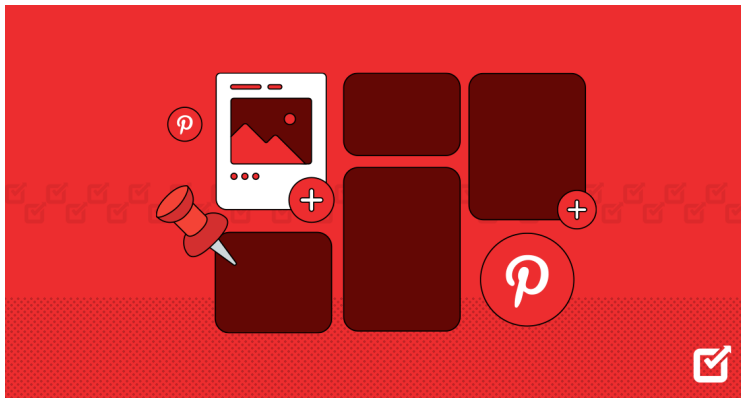
1 Machine Learning with Graphs

2 Introduction to Graphs

3 Graph Neural Networks

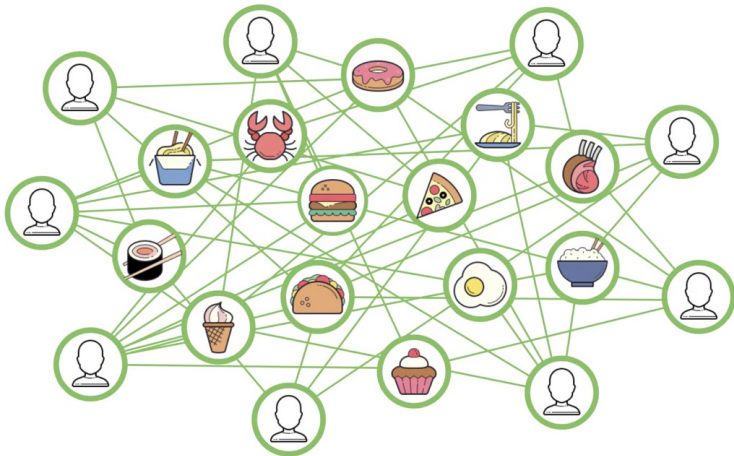
4 Training GNNs

# Recommendation System: Pinterest



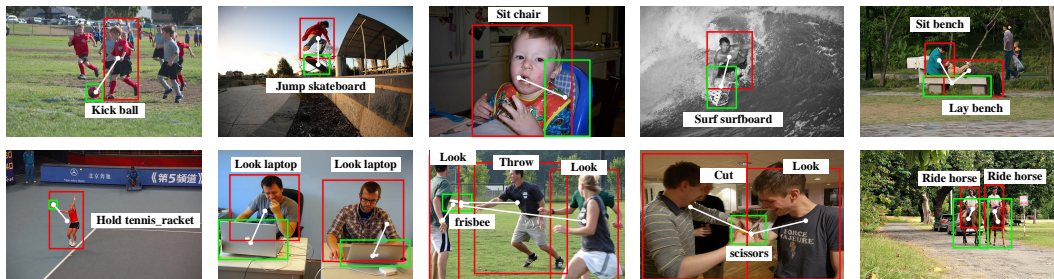
- Pinterest uses Graph Convolutional Networks (GCNs) to enhance recommendations by modeling user-item interactions as a bipartite graph, allowing the system to capture complex user preferences through aggregated neighbor information.

# Recommendation System: Uber Eats



- Uber Eats employs Graph Neural Networks (GNNs) to enhance its recommendation system by modeling the relationships between users, restaurants, and dishes as a graph, enabling the platform to provide more personalized and relevant food and restaurant recommendations to users.

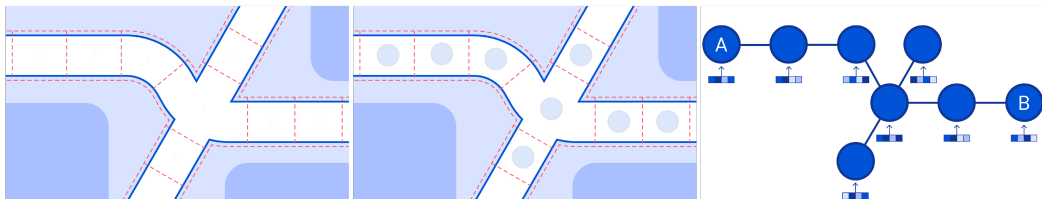
# Human-Object Interaction



- This paper introduces the Graph Parsing Neural Network (GPNN), which uses message passing to recognize human-object interactions in images and videos by dynamically constructing a parse graph that captures complex relationships between humans and objects.

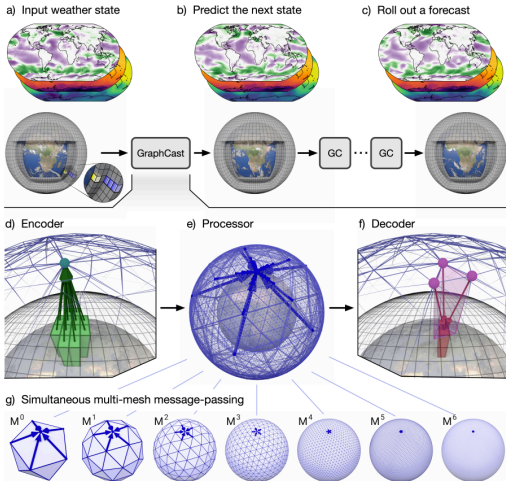


# Estimated Time of Arrival (ETA) Prediction with Graph Neural Networks



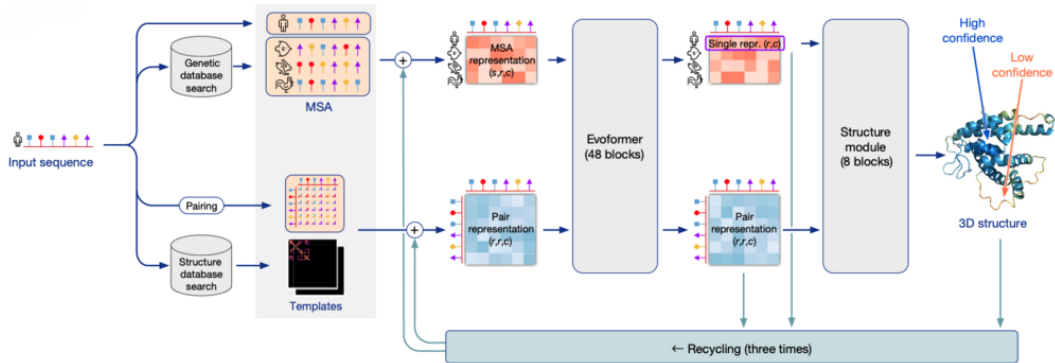
- Google Maps employs Graph Neural Networks (GNNs) to model road networks as graphs, using message passing to integrate spatial connections and real-time traffic data, enhancing ETA predictions on complex routes.

# GraphCast: Weather Forecasting



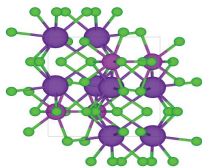
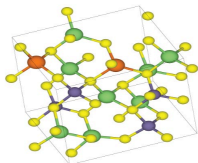
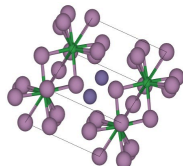
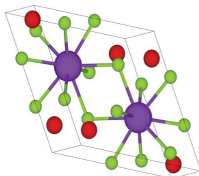
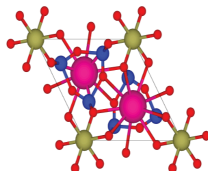
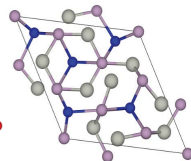
- Google DeepMind introduced GraphCast, a GNN-based weather model with an Encoder-Processor-Decoder architecture that captures global spatial dependencies for improved medium-range forecasts.

# Protein Folding with AlphaFold



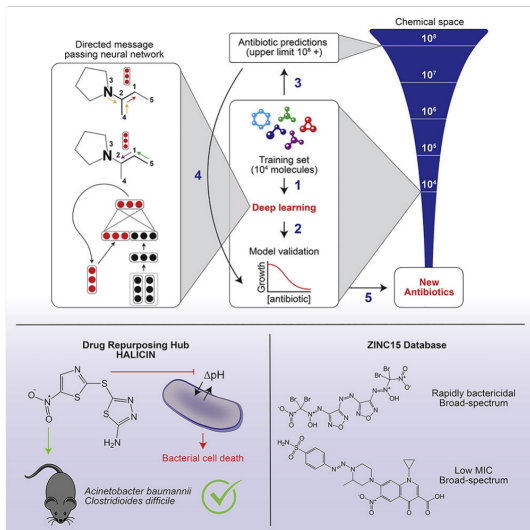
- AlphaFold 2 models relationships between amino acid residues using graph-inspired techniques, including attention mechanisms and pairwise representations. This approach captures essential relational information with customized attention layers optimized for protein folding.

## Materials Science

 $K_2BiCl_5$  $Li_4MgGe_2S_7$  $Mo_5GeB_2$  $KV_3Se_3$  $Rb_2HfSi_3O_9$  $Tm_5Pd_9P_7$ 

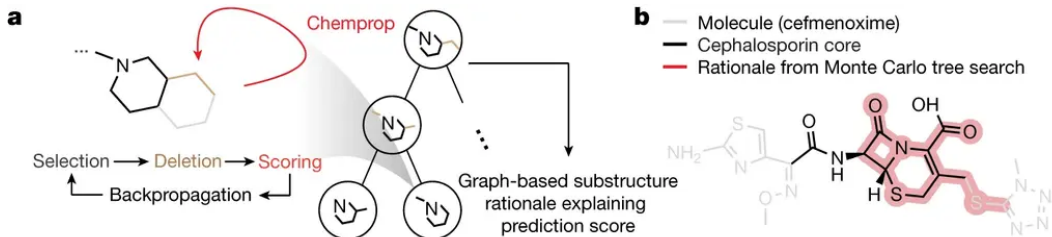
- DeepMind's AI tool, GNoME (Graph Network for Materials Exploration), uses graph neural networks (GNNs) to discover and predict the stability of new crystalline materials. This approach identified over 2 million new structures, with about 380,000 classified as stable.

# Drug Discovery: Graph Neural Networks for Antibiotic Discovery



- MIT's COLLINS LAB uses a Directed Message Passing Neural Network (D-MPNN) to model molecules as graphs (atoms as nodes, bonds as edges), enabling direct prediction of antibacterial efficacy and efficient exploration of chemical space.

## Explainable AI for Antibiotic Discovery



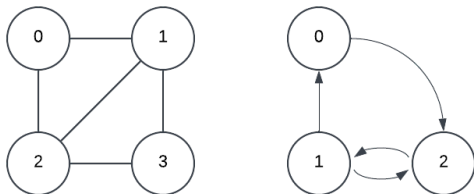
- Researchers from MIT and Harvard have developed an explainable AI system to identify antibiotic activity in molecules by combining a Graph Neural Network (GNN) with Monte Carlo Tree Search (MCTS). This approach isolates specific chemical substructures, known as *rationales*, that are associated with antibiotic effectiveness.

# Outline

- 1 Machine Learning with Graphs
- 2 Introduction to Graphs**
- 3 Graph Neural Networks
- 4 Training GNNs

# Graph Fundamentals

## Graph Definition:



- A graph  $\mathcal{G}$  is an ordered pair  $(\mathcal{V}, \mathcal{E})$ .
- $\mathcal{V}$  is a set of **vertices** (or nodes) representing entities like users or items.
- $\mathcal{E}$  is a set of **edges**, representing relationships between vertices, such as interactions.
- **Directed** edges are ordered pairs.
- **Undirected** edges are unordered pairs.

## Graph Representations:

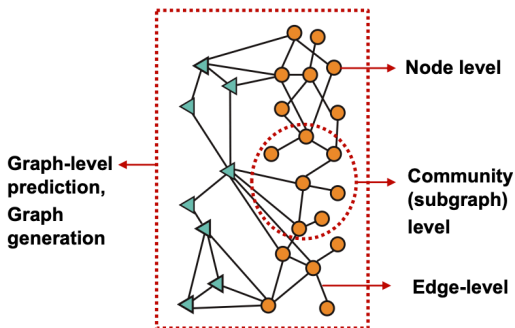
- **Adjacency Matrix:**  $A_{ij}$  indicates the presence or weight of an edge between nodes  $i$  and  $j$ .
- **Edge List:** A list of node pairs (or triplets, if edges have weights or attributes).

$$\mathbf{A}_1 = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 1 & 1 & 0 \end{bmatrix}, \quad E_2 = \{(0, 2), (1, 0), (1, 2), (2, 1)\}$$

where  $\mathbf{A}_1$  is the adjacency matrix for the first graph and  $E_2$  is the edge list for the second graph.



# Graph Learning Tasks



- **Node Level Tasks:**

- **Social Networks:** Predict user interests or groups.
- **Knowledge Graphs:** Classify entities in a knowledge base.

- **Edge Level Tasks:**

- **Recommendation Systems:** Predict user-item interactions (e.g., Uber Eats, Pinterest).
- **Social Networks:** Suggest new connections between users.

- **Graph Level Tasks:**

- **Drug Discovery:** Classify molecules based on properties (e.g., antibiotic activity).
- **Protein Function Prediction:** Predict biological function based on protein structures.

# Outline

- 1 Machine Learning with Graphs
- 2 Introduction to Graphs
- 3 Graph Neural Networks**
- 4 Training GNNs

# Using MLP on Graphs

## MLP for Node Feature Updates:

- Each node  $i$  in the graph has a **feature vector**  $\mathbf{x}_i \in \mathbb{R}^d$ .
- Node features can be updated by a Multi-Layer Perceptron (MLP):

$$\mathbf{h}_i = \phi(\mathbf{W}\mathbf{x}_i)$$

where  $\mathbf{W} \in \mathbb{R}^{m \times d}$  is the weight matrix,  $\phi$  is an activation function, and bias terms are omitted here for simplicity.

- In matrix form, the update becomes:

$$\mathbf{H} = \phi(\mathbf{X}\mathbf{W}^\top),$$

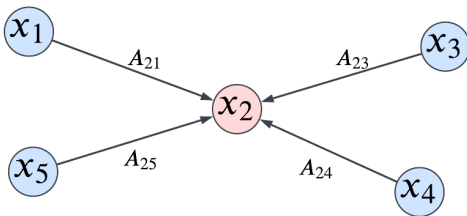
where:

$$\mathbf{X} = [\mathbf{x}_1 \quad \cdots \quad \mathbf{x}_n]^\top \in \mathbb{R}^{n \times d}, \quad \mathbf{H} = [\mathbf{h}_1 \quad \cdots \quad \mathbf{h}_n]^\top \in \mathbb{R}^{n \times m}.$$

## Limitation of Ignoring Graph Structure:

- This approach does not consider **graph structure**, treating each node as an isolated point.
- **Drawback:** Without adjacency information, we lose valuable relationships and structural insights of the graph, which are often crucial for graph-based tasks.

## Simple Aggregation with Neighbors



- Assume the graph is *unweighted* and *undirected*, i.e.,  $A_{ij} = A_{ji} = 1$  if edge  $(i, j)$  exists.
- The node feature is updated by **aggregating** with its neighbors:

$$h_i = \sigma \left( \sum_{j \in \mathcal{N}_i} \mathbf{W} \mathbf{x}_j \right) = \sigma \left( \sum_{j=1}^n A_{ij} \mathbf{W} \mathbf{x}_j \right),$$

where  $\mathcal{N}_i$  is the set of neighbors of node  $i$ .

- In matrix form, the update becomes:  $\mathbf{H} = \sigma(\mathbf{A}\mathbf{X}\mathbf{W}^\top)$
- A **deeper** neural network can be created by repeating this recurrent update:

$$\mathbf{H}^{(\ell)} = \sigma(\mathbf{A}\mathbf{H}^{(\ell-1)}\mathbf{W}^{(\ell)\top})$$

where  $\mathbf{H}^0 = \mathbf{X}$ .

# Graph Convolutional Network (GCN)

- Standard adjacency matrices omit **self-connections**, so we modify by adding self-loops:

$$\mathbf{h}_i^{(\ell)} = \sigma \left( \mathbf{W}^{(\ell)} \mathbf{h}_i^{(\ell-1)} + \sum_{j \in \mathcal{N}_i \setminus \{i\}} \mathbf{W}^{(\ell)} \mathbf{h}_j^{(\ell-1)} \right) \Rightarrow \mathbf{H}^{(\ell)} = \sigma \left( \hat{\mathbf{A}} \mathbf{H}^{(\ell-1)} \mathbf{W}^{(\ell)\top} \right)$$

where  $\hat{\mathbf{A}} = \mathbf{A} + \mathbf{I}$ .

- **Normalization** or averaging the aggregation prevents output scaling:

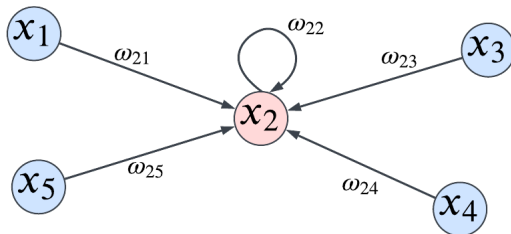
$$\mathbf{h}_i^{(\ell)} = \sigma \left( \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathbf{W}^{(\ell)} \mathbf{h}_j^{(\ell-1)} \right) \Rightarrow \mathbf{H}^{(\ell)} = \sigma \left( \hat{\mathbf{D}}^{-1} \hat{\mathbf{A}} \mathbf{H}^{(\ell-1)} \mathbf{W}^{(\ell)\top} \right)$$

where  $\hat{\mathbf{D}}_{ii} = \sum_{j=1}^n \hat{\mathbf{A}}_{ij}$  is the degree matrix.

- The **graph convolutional network (GCN)** apply **symmetric** normalization:

$$\mathbf{h}_i^{(\ell)} = \sigma \left( \sum_{j \in \mathcal{N}_i} \frac{1}{\sqrt{|\mathcal{N}_i|} \sqrt{|\mathcal{N}_j|}} \mathbf{W}^{(\ell)} \mathbf{h}_j^{(\ell-1)} \right) \Rightarrow \mathbf{H}^{(\ell)} = \sigma \left( \hat{\mathbf{D}}^{-1/2} \hat{\mathbf{A}} \hat{\mathbf{D}}^{-1/2} \mathbf{H}^{(\ell-1)} \mathbf{W}^{(\ell)\top} \right)$$

## GCN: Illustration



The node features are updated in **graph convolutional network (GCN)** using **symmetric** normalization:

$$\mathbf{h}_i^{(\ell)} = \sigma \left( \sum_{j \in \mathcal{N}_i} \omega_{ij} \mathbf{W}^{(\ell)} \mathbf{h}_j^{(\ell-1)} \right)$$

where  $\omega_{ij} := \hat{\mathbf{A}}_{ij} / \sqrt{\hat{\mathbf{D}}_{ii} \hat{\mathbf{D}}_{jj}}$ .

# Graph Attention Network (GAT)

GAT uses **attention** mechanisms to assign different levels of importance to neighboring nodes.

- Computes the **attention score** of node  $j$  to node  $i$  based on their features:

$$e_{ij}^{(\ell)} = a^{(\ell)}(\mathbf{h}_i^{(\ell-1)}, \mathbf{h}_j^{(\ell-1)})$$

where  $a$  is a *shared* attention mechanism that outputs attention scores for each edge.

- A common choice of  $a$  is a fully connected layer:

$$a^{(\ell)}(\mathbf{h}_i, \mathbf{h}_j) = \text{LeakyReLU}\left(\mathbf{a}^{(\ell)\top} [\mathbf{W}^{(\ell)} \mathbf{h}_i; \mathbf{W}^{(\ell)} \mathbf{h}_j]\right)$$

where  $[\cdot; \cdot]$  denotes concatenation, and  $\mathbf{a}^{(\ell)}$  and  $\mathbf{W}^{(\ell)}$  are learnable.

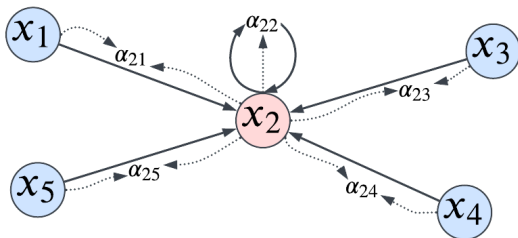
- Applies a softmax function to normalize attention scores

$$\alpha_i^{(\ell)} = \text{softmax}(e_i^{(\ell)})$$

- Aggregates **neighboring** features using normalized attention coefficients to update node features:

$$\mathbf{h}_i^{(\ell+1)} = \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^{(\ell)} \mathbf{W} \mathbf{h}_j^{(\ell)} \right)$$

## GAT: Illustration



- Computes the **attention score** of node  $j$  to node  $i$  based on their features:

$$e_{ij}^{(\ell)} = a^{(\ell)}(\mathbf{h}_i^{(\ell-1)}, \mathbf{h}_j^{(\ell-1)})$$

- Applies a softmax function to normalize attention scores

$$\alpha_i^{(\ell)} = \text{softmax}(e_i^{(\ell)})$$

- Aggregates **neighboring** features using normalized attention coefficients to update node features:

$$\mathbf{h}_i^{(\ell+1)} = \sigma \left( \sum_{j \in \mathcal{N}_i} \alpha_{ij}^{(\ell)} \mathbf{W} \mathbf{h}_j^{(\ell)} \right)$$



# Message Passing Neural Network (MPNN)

MPNNs extend GCNs and GATs by incorporating **edge features** (e.g., chemical bonds).

- Given node features  $\mathbf{h}_i^{(l-1)}$ ,  $\mathbf{h}_j^{(l-1)}$ , and edge features  $\mathbf{e}_{ij}$ , the **message function**  $f_e^{(\ell)}$  computes:

$$\mathbf{m}_{ij}^{(l)} = f_e^{(l)}(\mathbf{h}_i^{(l-1)}, \mathbf{h}_j^{(l-1)}, \mathbf{e}_{ij}).$$

- A common choice for  $f_e$  is a fully connected layer:

$$f_e^{(l)}(\mathbf{h}_i^{(l-1)}, \mathbf{h}_j^{(l-1)}, \mathbf{e}_{ij}) = \sigma \left( \mathbf{W}_e^{(l)} [\mathbf{h}_i^{(l-1)}; \mathbf{h}_j^{(l-1)}; \mathbf{e}_{ij}] \right)$$

where  $[\cdot; \cdot]$  denotes concatenation and  $\mathbf{W}_e^{(\ell)}$  is learnable matrix.

- The node feature  $\mathbf{h}_i^{(l)}$  is updated by **aggregating messages** from neighbors:

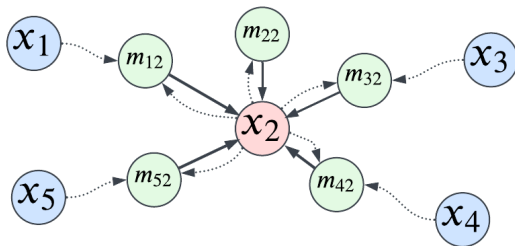
$$\mathbf{h}_i^{(l)} = f_v^{(l)}(\mathbf{h}_i^{(l-1)}, \{\mathbf{m}_{ij}^{(l)}\}_j)$$

- A common aggregation function  $f_v$  sums messages:

$$f_v^{(l)}(\mathbf{h}_i^{(l-1)}, \{\mathbf{m}_{ij}^{(l)}\}_j) = \sigma \left( \mathbf{W}_v^{(l)} \mathbf{h}_i^{(l-1)} + \mathbf{W}_m^{(l)} \sum_{j \in \mathcal{N}_i} \mathbf{m}_{ij}^{(l)} \right)$$

where  $\mathbf{W}_v^{(l)}$  and  $\mathbf{W}_m^{(l)}$  are learnable matrices.

## MPNN: Illustration



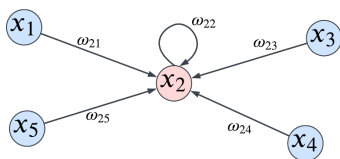
- Given node features  $\mathbf{h}_i^{(l-1)}$ ,  $\mathbf{h}_j^{(l-1)}$ , and edge features  $e_{ij}$ , the **message function**  $f_e^{(\ell)}$  computes:

$$\mathbf{m}_{ij}^{(l)} = f_e^{(l)}(\mathbf{h}_i^{(l-1)}, \mathbf{h}_j^{(l-1)}, e_{ij}).$$

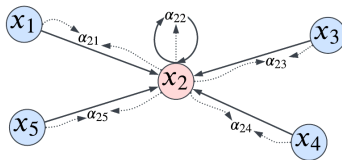
- The node feature  $\mathbf{h}_i^{(l)}$  is updated by **aggregating messages** from neighbors:

$$\mathbf{h}_i^{(l)} = f_v^{(l)}(\mathbf{h}_i^{(l-1)}, \{\mathbf{m}_{ij}^{(l)}\}_j)$$

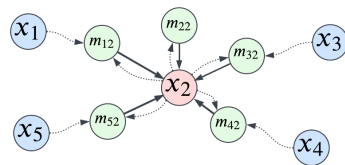
## Summary



GCNs



GATs



MPNNs

# Outline

- 1 Machine Learning with Graphs
- 2 Introduction to Graphs
- 3 Graph Neural Networks
- 4 Training GNNs**

# Node Classification with GNNs

**Goal:** Predict a label for each node in the graph.

- **Node Representation:** After  $L$  layers, each node  $i$  has an updated representation:

$$\mathbf{h}_i^{(L)} = f_{\text{GNN}}(\mathbf{X}, \mathbf{A})$$

where  $\mathbf{X}$  is the node feature matrix and  $\mathbf{A}$  is the adjacency matrix.

- **Prediction:** Apply a classifier to the final node representation:

$$\hat{y}_i = \text{softmax}(\mathbf{W}_{\text{out}} \mathbf{h}_i^{(L)})$$

where  $\mathbf{W}_{\text{out}}$  is a learnable weight matrix for classification.

- **Loss Function:** Use cross-entropy loss to compare predictions with ground truth labels:

$$\mathcal{L} = - \sum_{i \in \mathcal{V}} y_i \log(\hat{y}_i)$$

where  $y_i$  is the true label of node  $i$ .

# Graph Classification with GNNs

**Goal:** Predict a label for the entire graph.

- **Node Representations:** After  $L$  layers, each node  $i$  has a final representation:

$$\mathbf{h}_i^{(L)} = f_{\text{GNN}}(\mathbf{X}, \mathbf{A})$$

- **Graph-Level Representation:** Aggregate all node representations to form a single graph vector:

$$\mathbf{h}_{\text{graph}} = \text{AGGREGATE}(\{\mathbf{h}_i^{(L)} | i \in \mathcal{V}\})$$

Common choices for AGGREGATE include mean, sum, or max pooling.

- **Prediction:** Apply a classifier to the aggregated graph representation:

$$\hat{y}_{\text{graph}} = \text{softmax}(\mathbf{W}_{\text{out}} \mathbf{h}_{\text{graph}})$$

- **Loss Function:** Use cross-entropy loss to match predictions with true graph labels:

$$\mathcal{L} = -y_{\text{graph}} \log(\hat{y}_{\text{graph}})$$

where  $y_{\text{graph}}$  is the true graph label.

# Link Prediction with GNNs

**Goal:** Predict the existence of an edge between two nodes.

- **Node Representations:** After  $L$  layers, each node  $i$  has a final representation:

$$\mathbf{h}_i^{(L)} = f_{\text{GNN}}(\mathbf{X}, \mathbf{A})$$

- **Edge Representation:** For a pair of nodes  $i$  and  $j$ , compute a combined feature vector to represent the link:

$$\mathbf{z}_{ij} = g(\mathbf{h}_i^{(L)}, \mathbf{h}_j^{(L)})$$

where  $g$  is a function such as concatenation  $[\mathbf{h}_i^{(L)}; \mathbf{h}_j^{(L)}]$ , element-wise product  $\mathbf{h}_i^{(L)} \odot \mathbf{h}_j^{(L)}$ , or distance-based functions.

- **Prediction:** Apply a scoring function to predict the likelihood of a link:

$$\hat{y}_{ij} = \sigma(\mathbf{W}_{\text{out}} \mathbf{z}_{ij})$$

where  $\sigma$  is the sigmoid function.

- **Loss Function:** Use binary cross-entropy to compare predictions with actual labels:

$$\mathcal{L} = - \sum_{(i,j) \in \mathcal{E}} y_{ij} \log(\hat{y}_{ij}) + (1 - y_{ij}) \log(1 - \hat{y}_{ij})$$

where  $y_{ij}$  indicates if a link exists between  $i$  and  $j$ .