

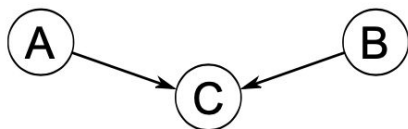
CSE8803/CX4803

Machine Learning in Computational Biology

Lecture 17:
Deep Learning for Networks
(Graph Neural Networks)

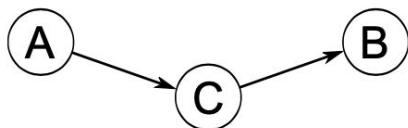
Yunan Luo

Conditional independence in BN



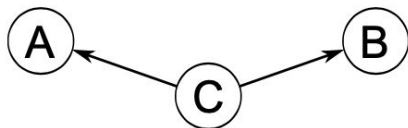
Case 1:

$$\begin{aligned}P(A, B, C) &= P(C|A, B)P(A, B) \\ &= P(C|A, B)P(A)P(B)\end{aligned}$$



Case 2:

$$\begin{aligned}P(A, B, C) &= P(B, C|A)P(A) \\ &= P(B|A, C)P(C|A)P(A) \\ &= P(B|C)P(C|A)P(A)\end{aligned}$$



Case 3:

$$\begin{aligned}P(A, B, C) &= P(A, B|C)P(C) \\ &= P(A|C)P(B|C)P(C)\end{aligned}$$

Example: proof of conditional independence of A, B in case 3

$$P(A, B|C) = \frac{P(A, B, C)}{P(C)} = \frac{P(C)P(A|C)P(B|C)}{P(C)} = P(A|C)P(B|C)$$

Conditional independence in BN

Each node is dependent only on its parents

Define:

$P(\mathbf{X})$ - the global probability of all variables

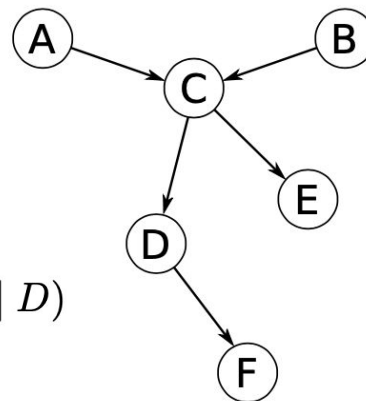
$P(\mathbf{X}) = P(A, B, C, D, E, F)$

$P(\mathbf{X}) = P(A) P(B) P(C | A, B) P(D | C) P(E | C) P(F | D)$

In general:

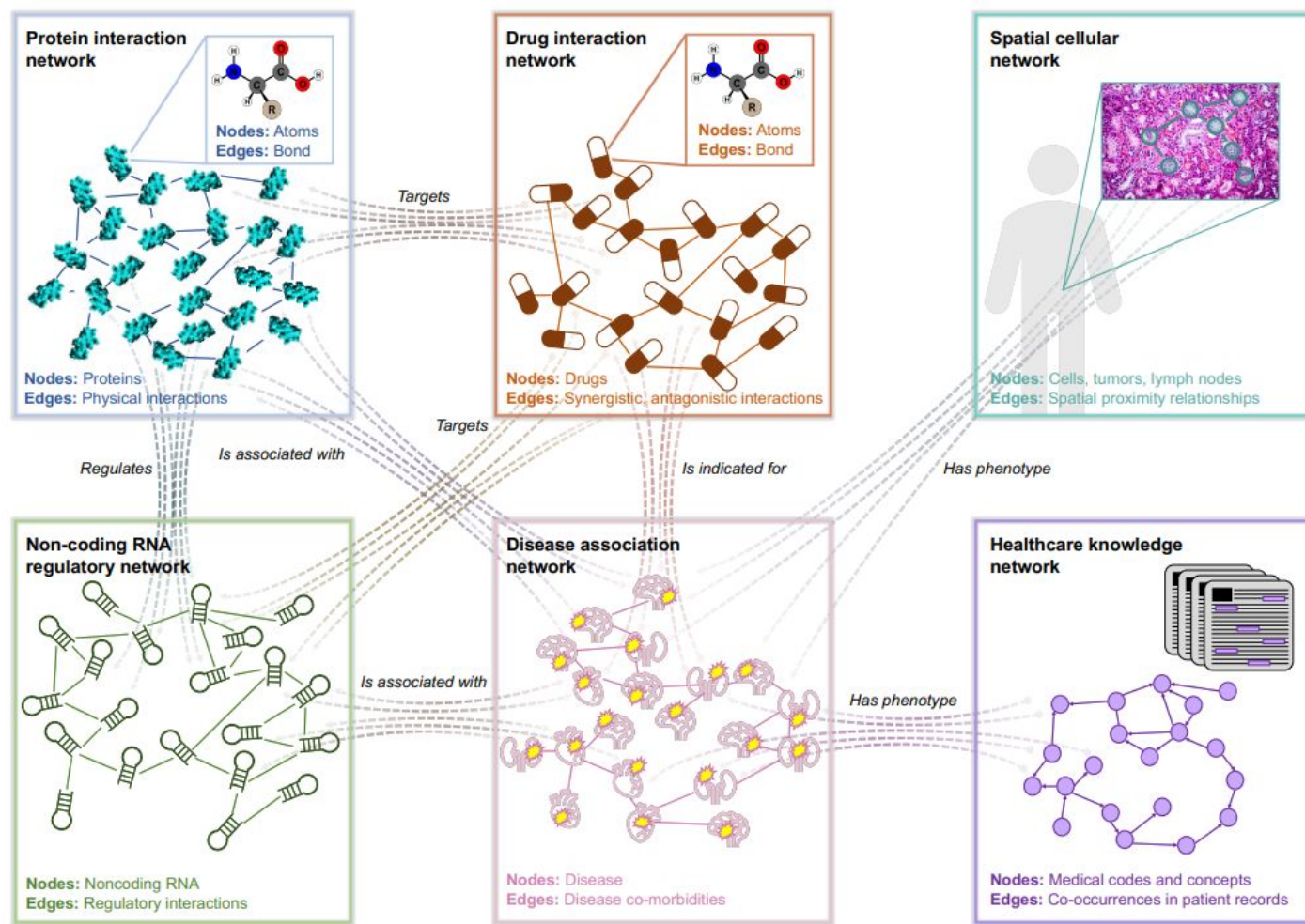
$$p(X_1, \dots, X_n) = \prod_{i=1}^n p(X_i | X_{\text{pa}(i)})$$

where $\text{pa}(i)$ is the parents of node i

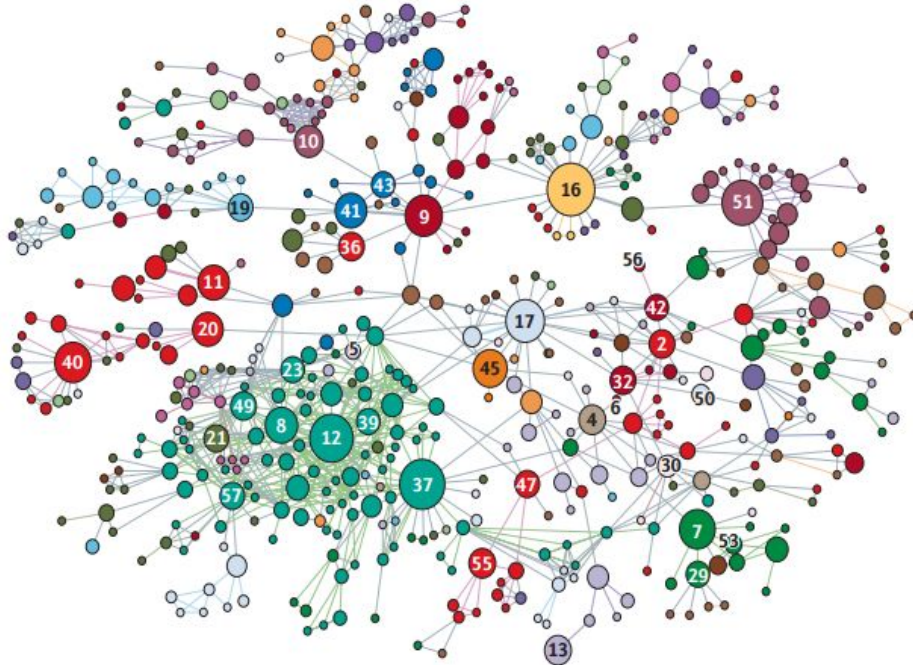


Each variable is conditionally independent of its non-descendants, given its parents.

Why networks in biology?



Human disease network



Node: protein
Edge: protein-protein interaction

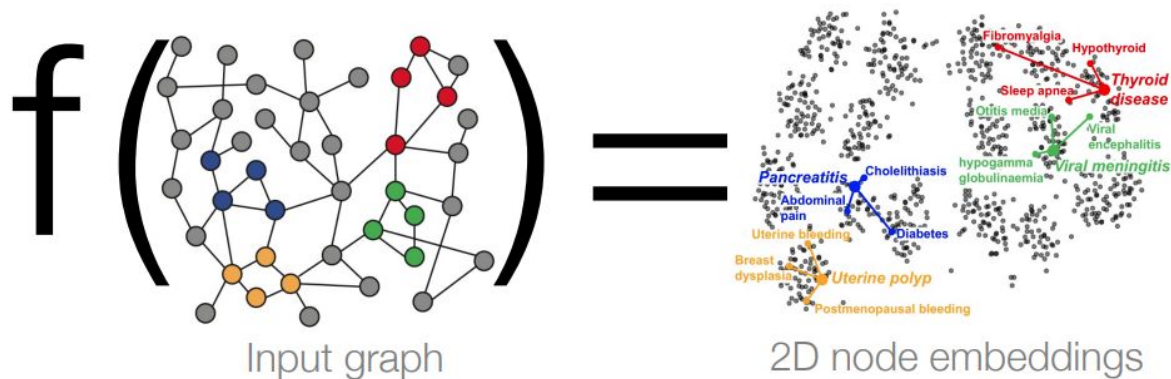
① Aldosteronism	②① Epilepsy	④② Myocardial infarction
② Alzheimer's disease	②② Fanconi's anaemia	④③ Myopathy
③ Anaemia, congenital deserythropoietic	②③ Fatty liver	④④ Nucleoside phosphorylase deficiency
④ Asthma	②④ Gastric cancer	④⑤ Obesity
⑤ Ataxia-telangiectasia	②⑤ Gilbert's syndrome	④⑥ Paraganglioma
⑥ Atherosclerosis	②⑥ Glaucoma 1A	④⑦ Parkinson's disease
⑦ Blood group	②⑦ Goitre congenital	④⑧ Pheochromocytoma
⑧ Breast cancer	②⑧ HARP syndrome	④⑨ Prostate cancer
⑨ Cardiomyopathy	②⑨ HELLP syndrome	④⑩ Pseudohypoaldosteronism
⑩ Cataract	②⑩ Haemolytic anaemia	④⑪ Retinitis pigmentosa
⑪ Charcot-Marie-Tooth disease	②⑪ Hirschprung disease	④⑫ Schizoaffective disorder
⑫ Colon cancer	②⑫ Hyperbilirubinaemia	④⑬ Spherocytosis
⑬ Complement component deficiency	②⑬ Hypertension	④⑭ Spina bifida
⑭ Coronary artery disease	②⑭ Hypertension diastolic	④⑮ Spinocerebellar ataxia
⑮ Coronary spasm	②⑮ Hyperthyroidism	④⑯ Stroke
⑯ Deafness	②⑯ Hypoaldosteronism	④⑰ Thyroid carcinoma
⑰ Diabetes mellitus	②⑰ Leigh syndrome	④⑱ Total iodide organification defect
⑱ Enolase-β deficiency	②⑱ Leukaemia	④⑲ Trifunctional protein deficiency
⑲ Epidermolysis bullosa	②⑲ Low renin hypertension	④⑳ Unipolar depression
	③① Lymphoma	
	③② Mental retardation	
	③③ Muscular dystrophy	

Proteins involved in the same disease have an increased tendency to interact with each other

Barabási et al, "Network medicine: a network-based approach to human disease", 2011

Recap: Network embeddings

- **Idea:** Map nodes to d-dimensional embeddings such that similar nodes in the graph are embedded close together

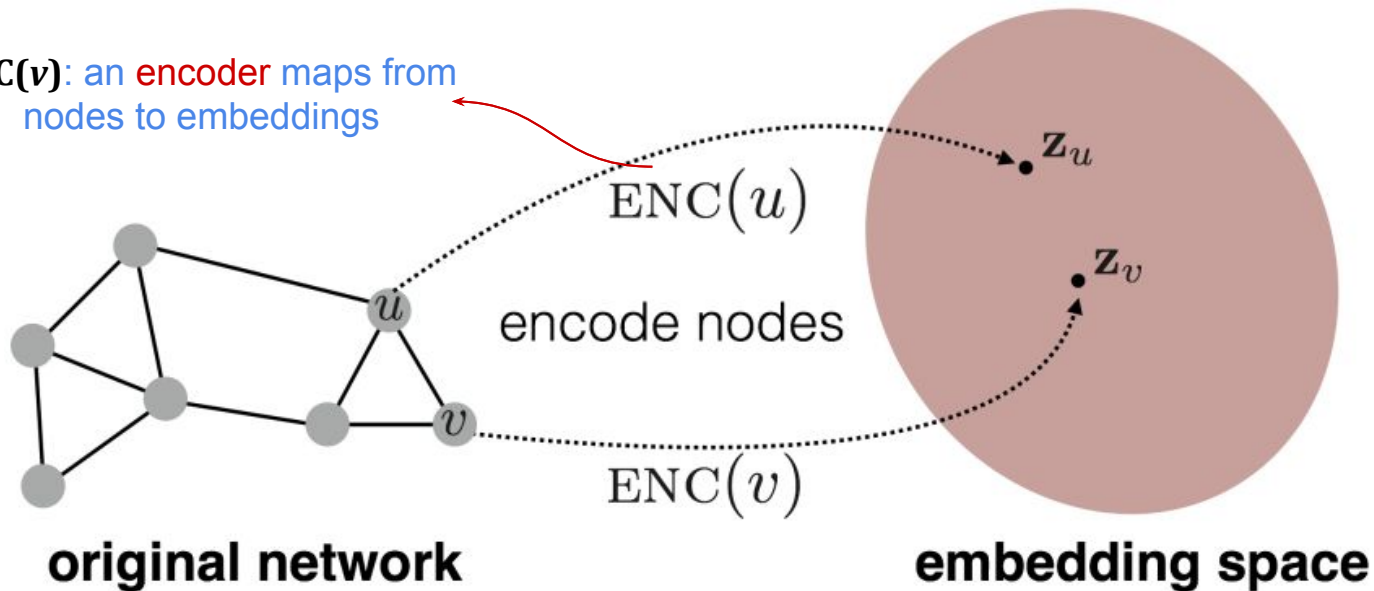


How to learn mapping function f ?

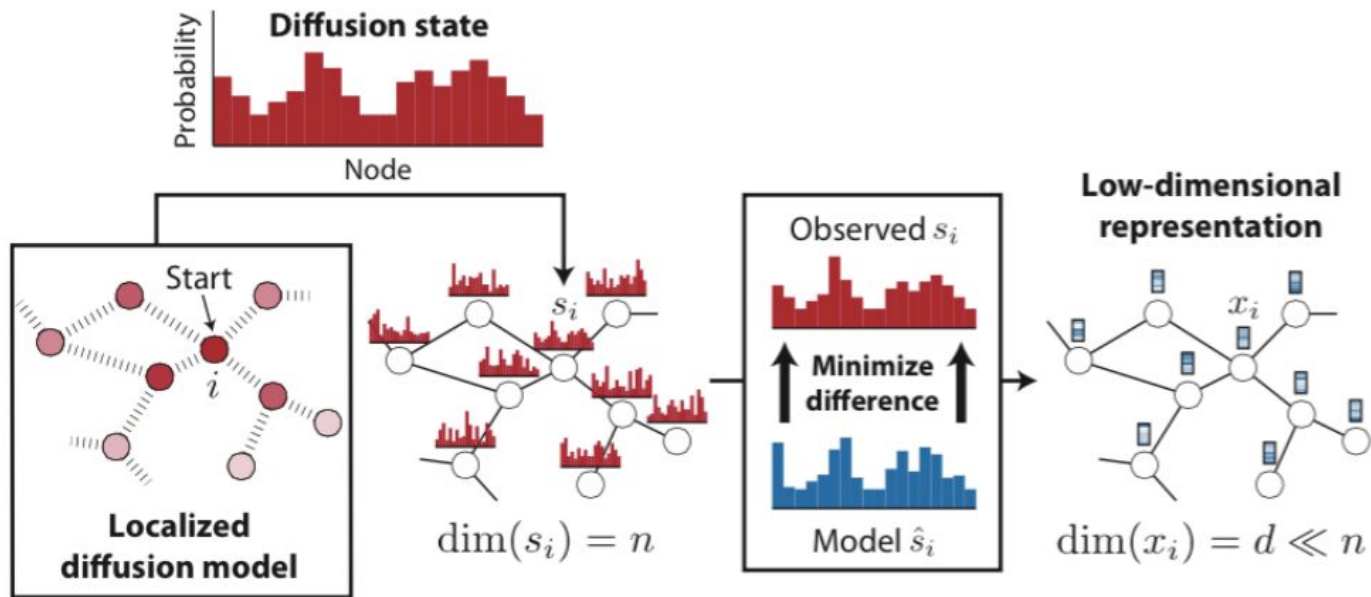
Recap: Embedding nodes

$$\underset{\text{In the original network}}{\text{similarity}(u, v)} \approx \underset{\text{Similarity of the embeddings}}{\mathbf{z}_v^T \mathbf{z}_u}$$

$\text{ENC}(v)$: an **encoder** maps from
nodes to embeddings

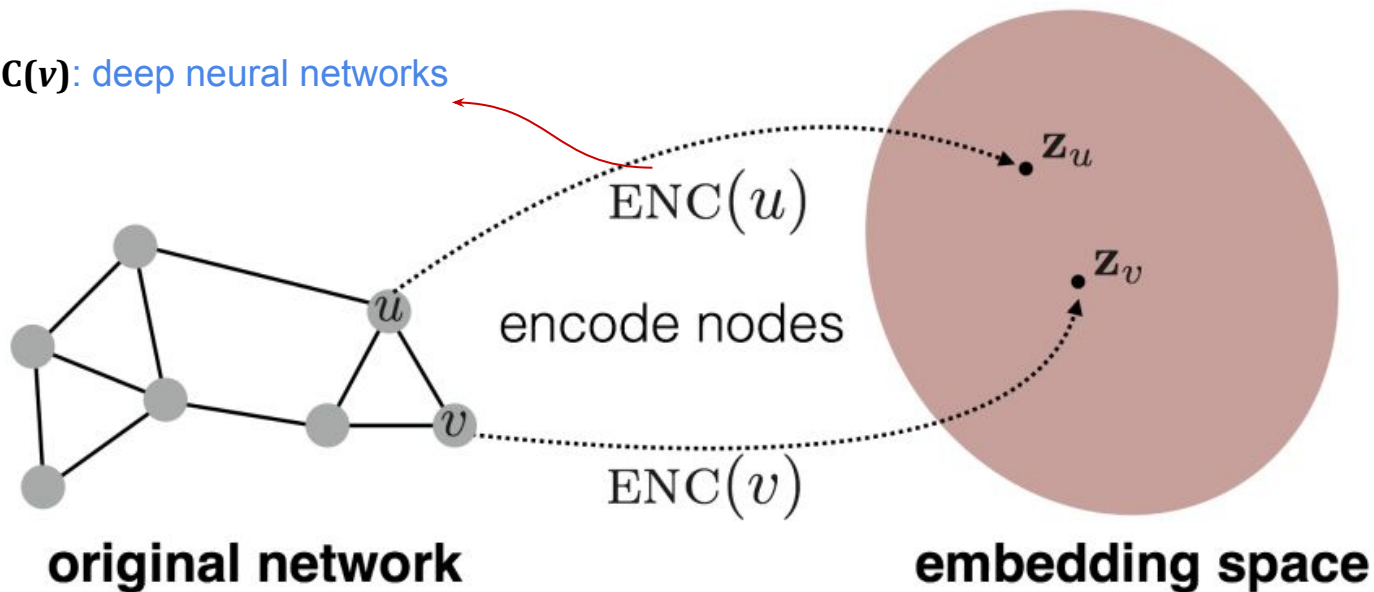


Recap: Diffusion-based approaches

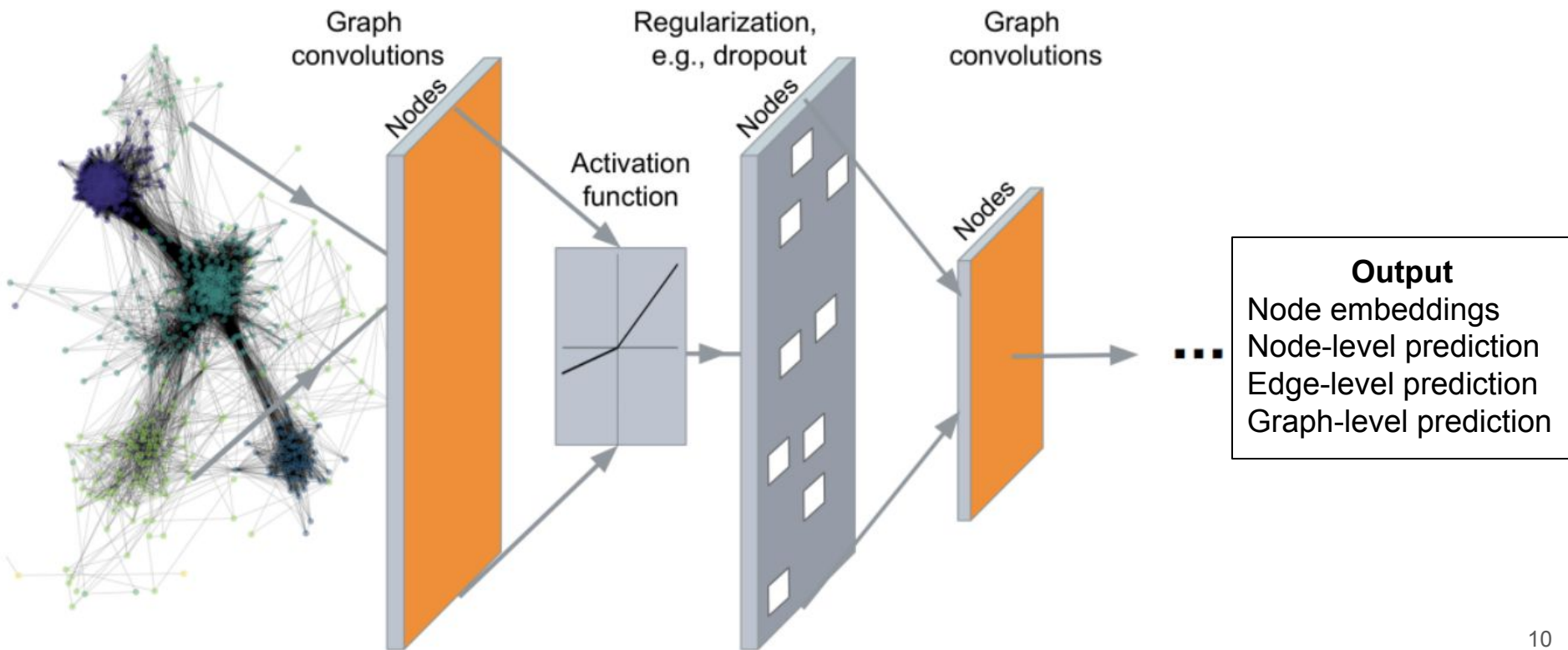


Today: Deep learning for graphs

$\text{ENC}(v)$: deep neural networks



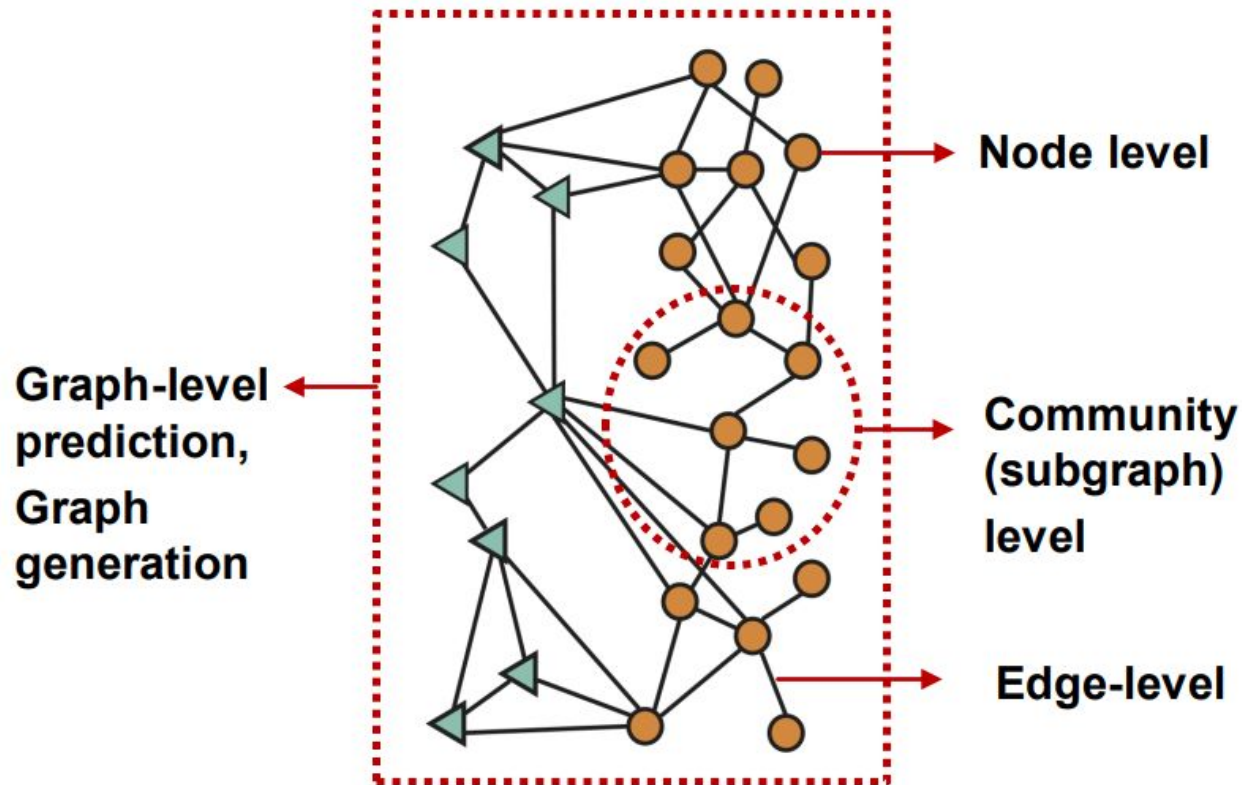
Graph Neural Network (GNN)



Problem setup

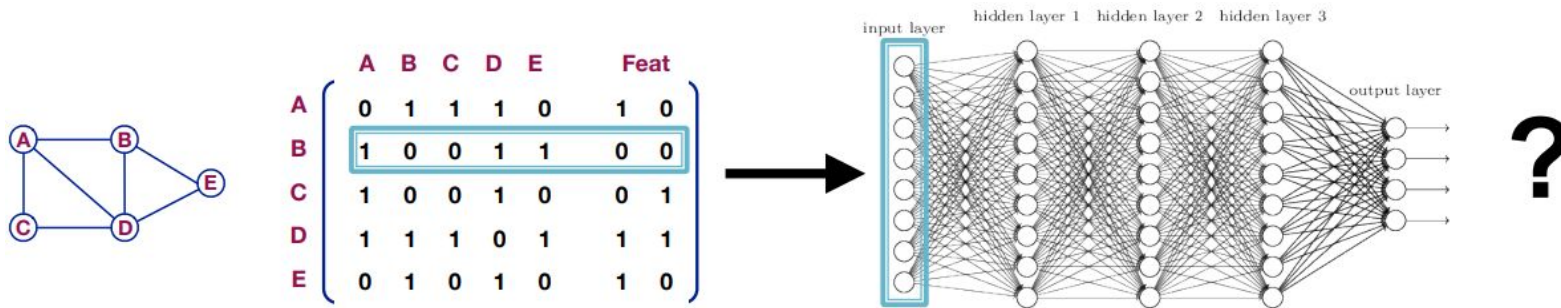
- Given a graph \mathbf{G}
 - \mathbf{V} is the **node set**
 - \mathbf{A} is the **adjacency matrix** (assume binary)
 - $\mathbf{X} \in \mathbb{R}^{m \times |V|}$ is a matrix of **node features**
 - v : a node in \mathbf{V}
 - $N(v)$: the set of neighbors of v
- Node features:
 - Social networks: User profile, User image
 - Biological networks: Gene expression profiles, gene functional information
 - When there is no node feature in the graph dataset:
 - Indicator vectors (one-hot encoding of a node)
 - Vector of constant 1: $[1, 1, \dots, 1]$

Recap: Different types of ML tasks on graphs



Idea 1: Fully-connected neural network

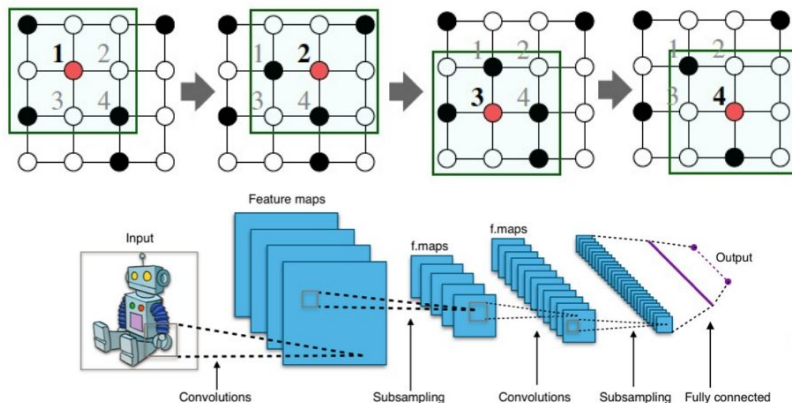
- Concatenate adjacency matrix and features
- Feed them into a deep neural network



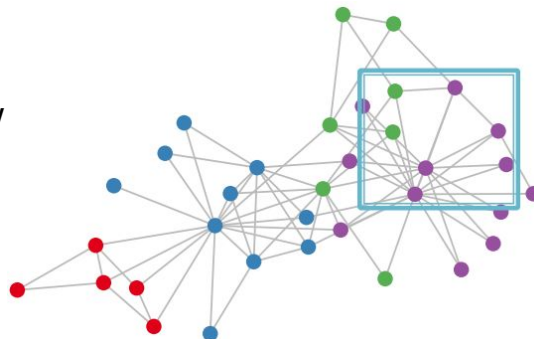
- Issues with this idea?
 - $O(|V|)$ parameters
 - Not applicable to graphs of different sizes
 - Sensitive to node ordering

Idea 2: Convolutional neural network

- CNN for image data

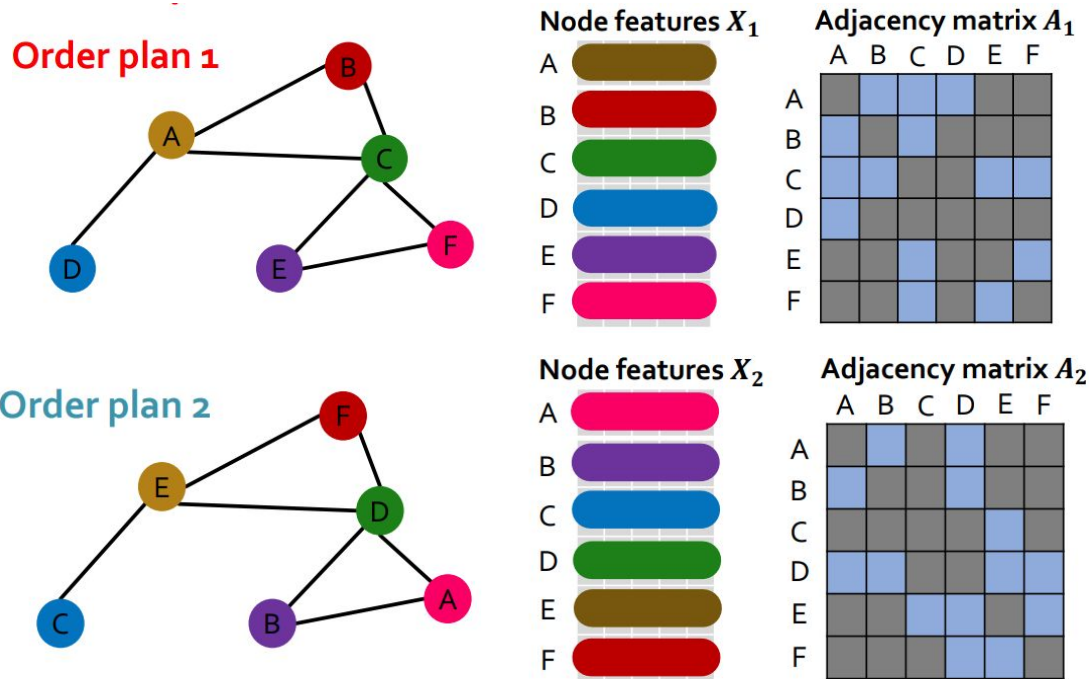


- CNN for Graph data?
 - No fixed notion of locality or sliding window on the graph
 - Graph is permutation invariant



Permutation invariance

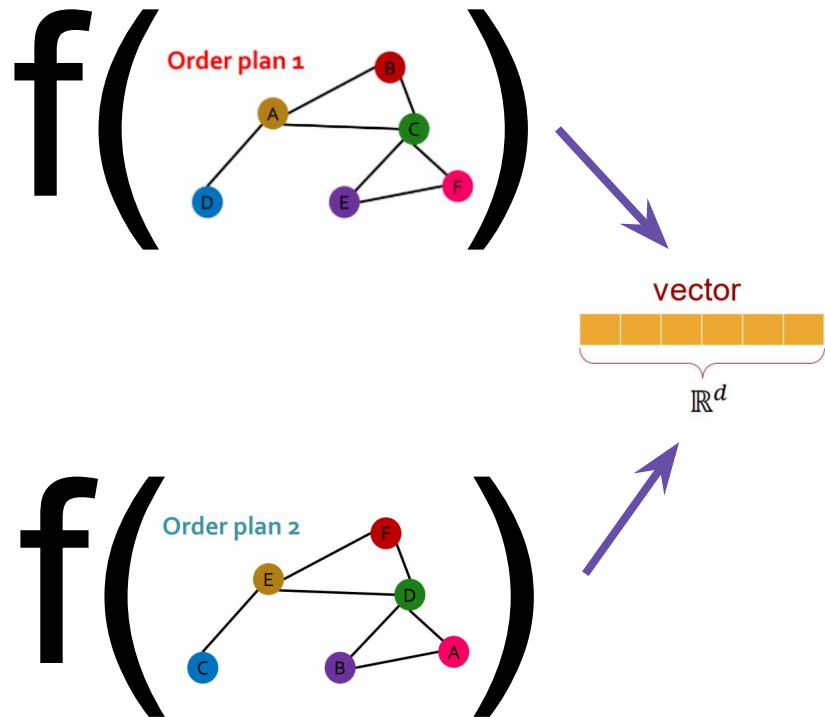
- Graph does not have a canonical order of the nodes



Graph/node representations should be the same for **Order plan 1** and **Order plan 2**

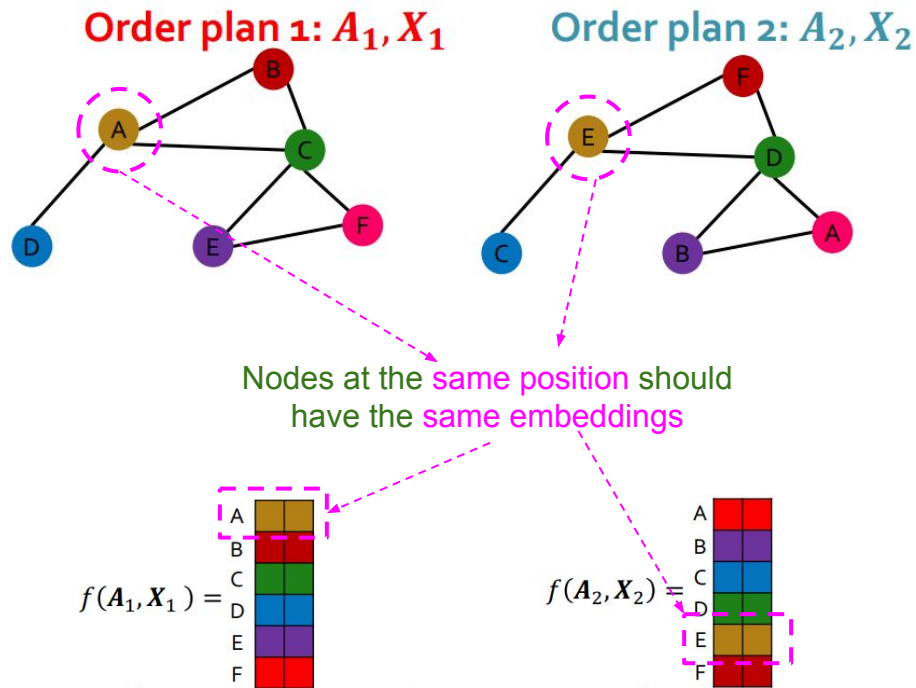
Permutation invariance

- Consider **graph** representation
- Goal**: learn a function f that maps a graph $G = (A, X)$ to a vector R^d
- If $f(A_i, X_i) = f(A_j, X_j)$ for any order plan i and j , we say f is a **permutation invariant function**



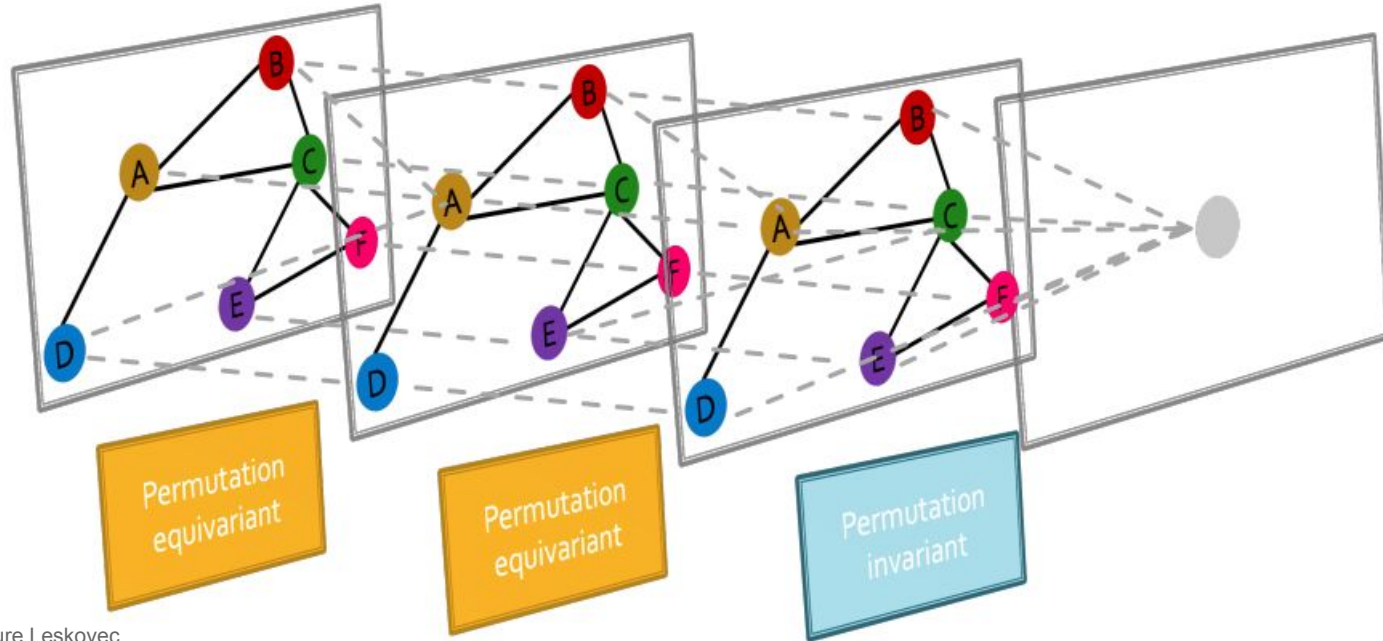
Permutation equivariance

- Consider **node** representation
- Goal learn a function f that maps a graph $G = (A, X)$ to a vector $R^{m \times d}$
 - m : #nodes, each row is the embedding of a node
- If every pair of nodes at the same position have the same embedding, we say f is a **permutation equivariant function**

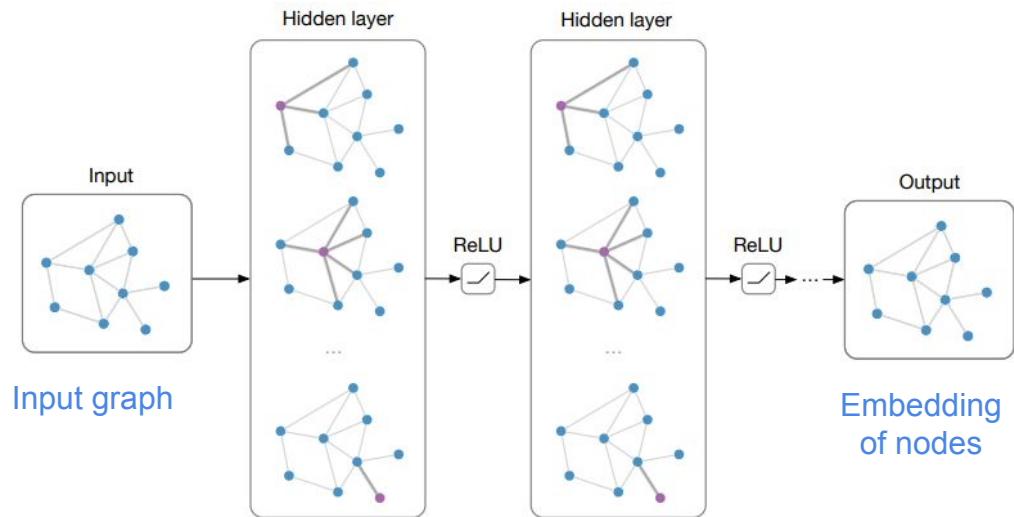


Graph neural networks overview

- Graph neural networks consist of multiple permutation equivariant / invariant functions



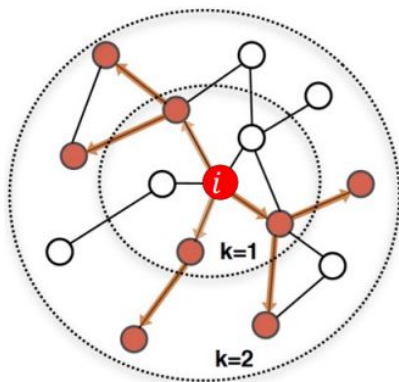
Graph neural network



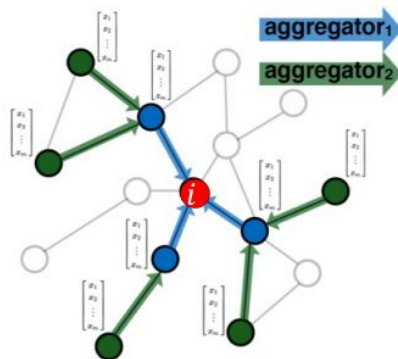
- **Main Idea:** Pass messages between pairs of nodes and agglomerate
- **Alternative Interpretation:** Pass messages between nodes to refine node (and possibly edge) representations

Graph neural network

Idea: Node's neighborhood defines a **computation graph**



Determine node
computation graph

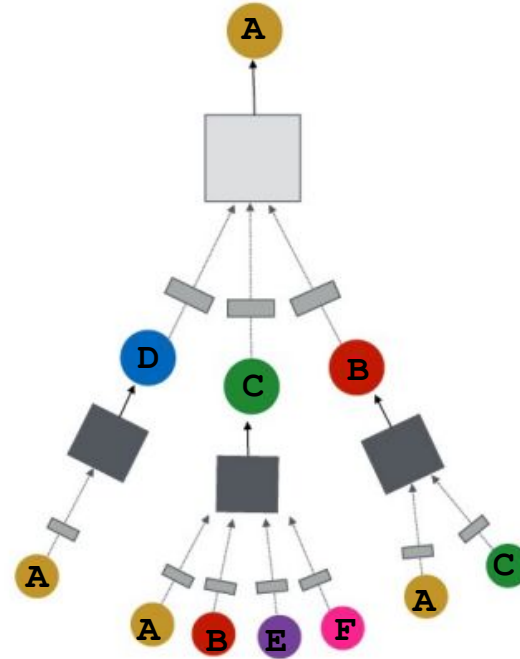
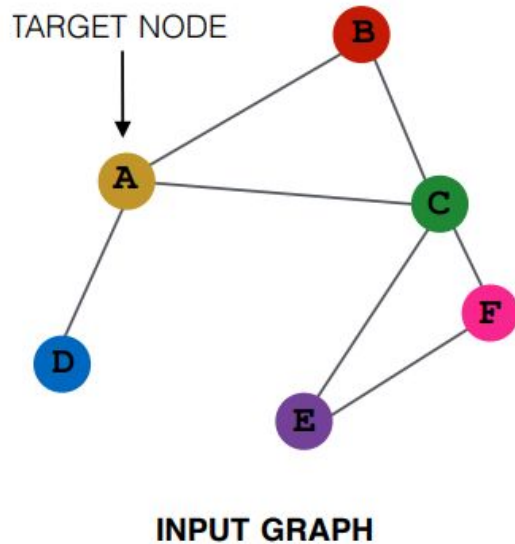


Propagate and
transform information

Goal: Learn how to propagate information across the graph to compute node features

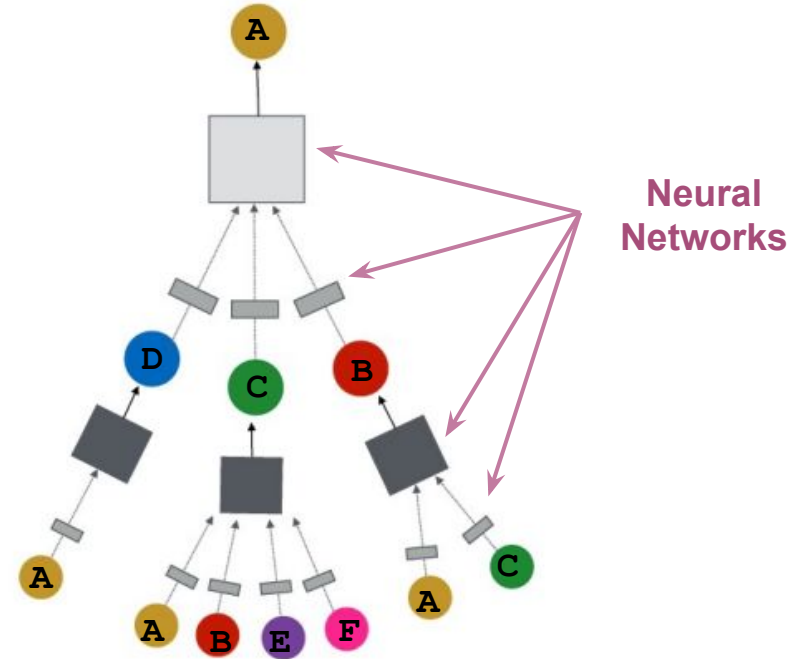
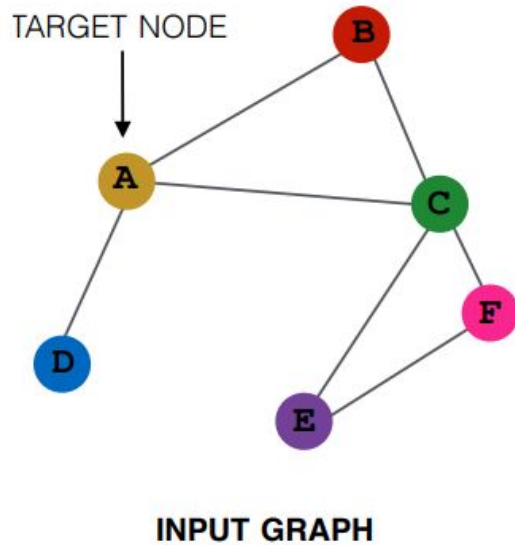
Idea: aggregate information from neighbors

- **Key idea:** generate node embeddings based on **local network neighborhoods**



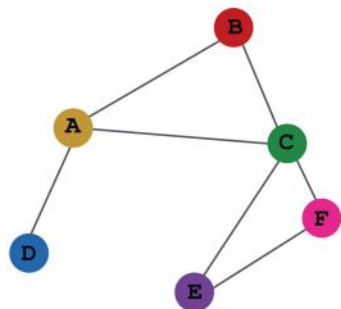
Idea: aggregate information from neighbors

- **Intuition:** Nodes aggregate information from their neighbors using neural networks



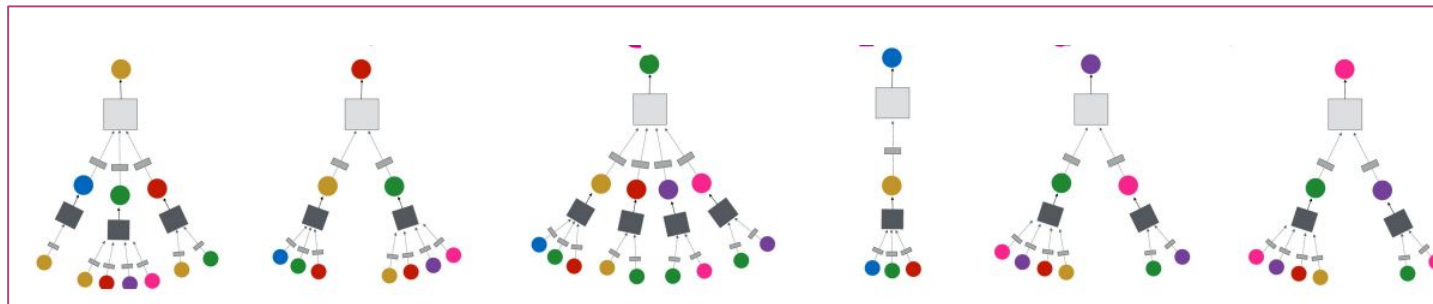
Idea: aggregate information from neighbors

Intuition: Node's neighborhood defines a **computation graph**

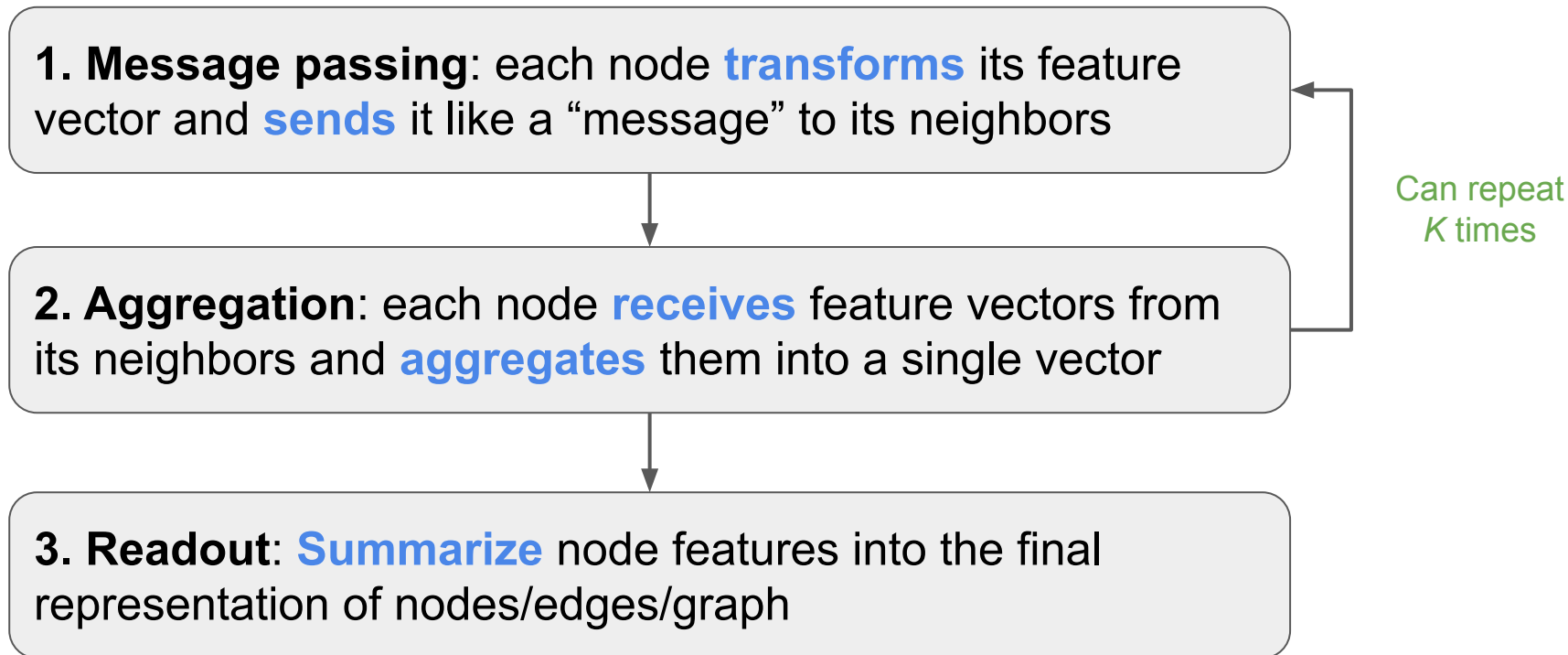


INPUT GRAPH

Every node defines a computation graph based on its neighborhood!

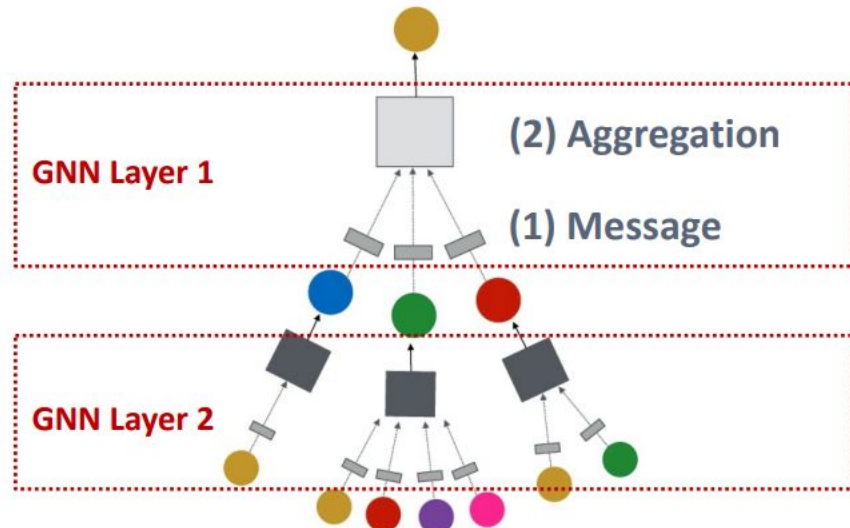


General GNN framework

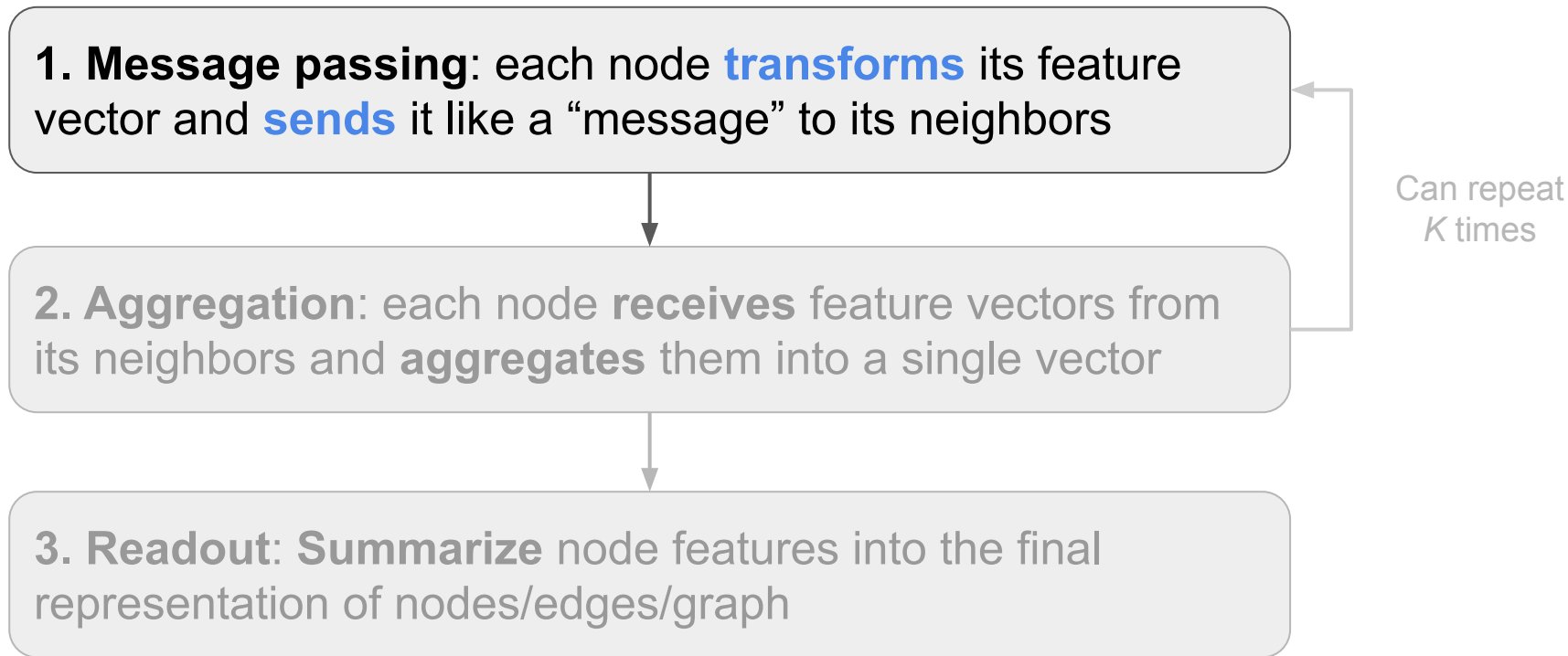


GNN can have many layers

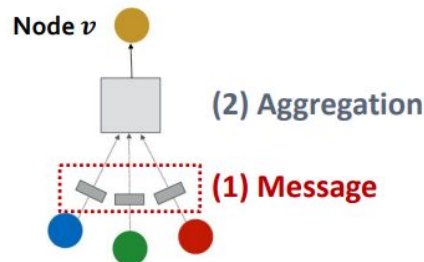
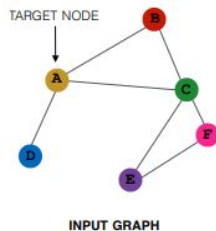
- GNN can be of arbitrary depth
- Nodes have embeddings at each layer
- Layer-0 embedding of node v is its input feature, \mathbf{x}_v
- Layer- k embedding gets information from nodes that are k hops away



General GNN framework



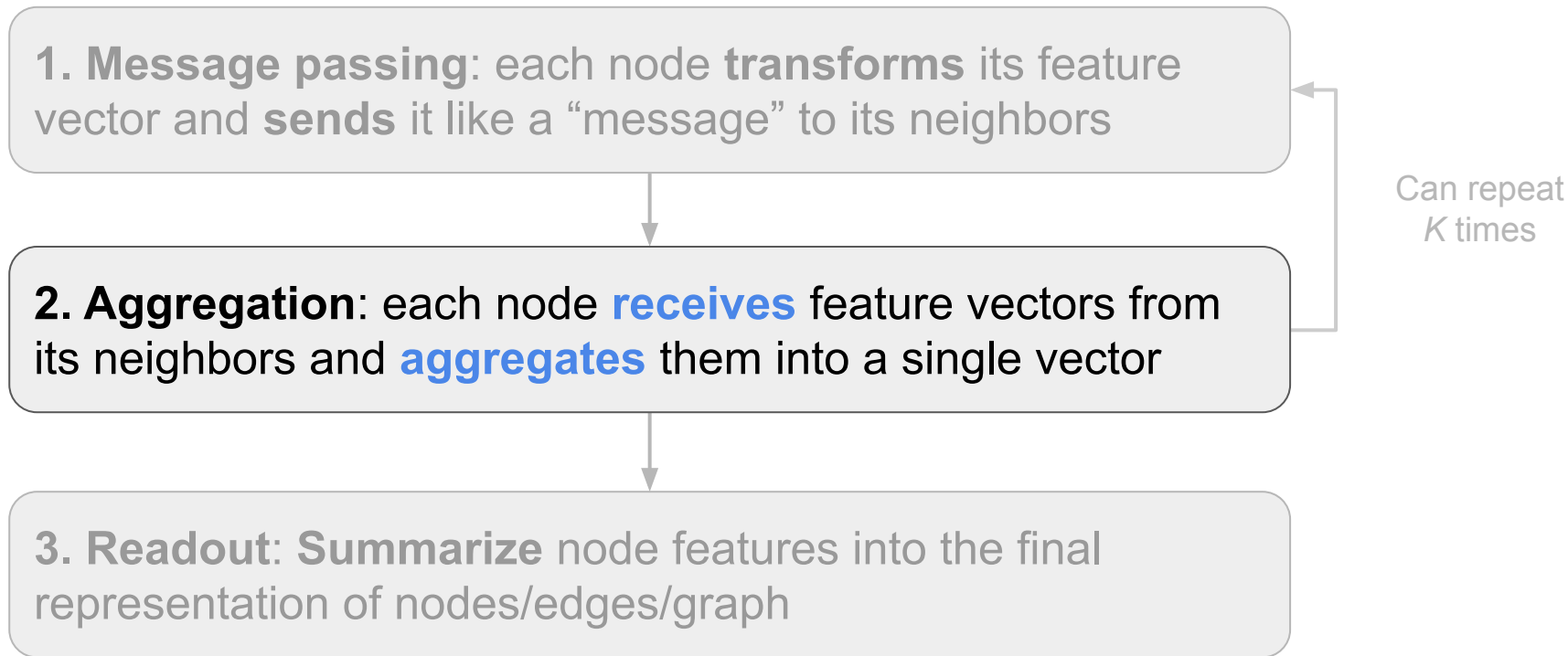
Message passing



The l -th layer of GNN

- **Message function:** $\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}(\mathbf{h}_u^{(l-1)})$
 - **Intuition:** Each node will create a message, which will be sent to other nodes later
 - **Example:** A Linear layer $\mathbf{m}_u^{(l)} = \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}$
 - Multiply node features with weight matrix $\mathbf{W}^{(l)}$
- ❖ At the 0-th step, \mathbf{h}_u^0 is simply the node feature \mathbf{x}_u

General GNN framework



Aggregation

- Aggregation: each node **receives** feature vectors from its neighbors and **aggregates** them into a single vector

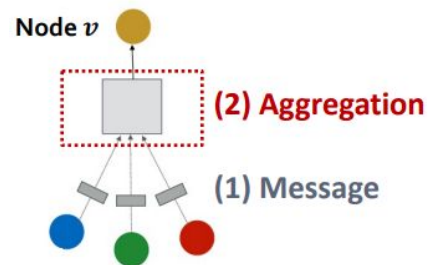
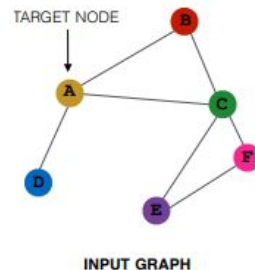
- **Intuition:** Each node will aggregate the messages from node v 's neighbors

$$\mathbf{h}_v^{(l)} = \text{AGG}^{(l)} \left(\left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\} \right)$$

aggregation
function

- **Example:** Sum(\cdot), Mean(\cdot) or Max(\cdot) aggregator

- $\mathbf{h}_v^{(l)} = \text{AGG}^{(l)} \left(\left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\}, \mathbf{m}_v^{(l)} \right)$



A single GNN layer

Summary:

- **(1) Message:** each node computes its own message

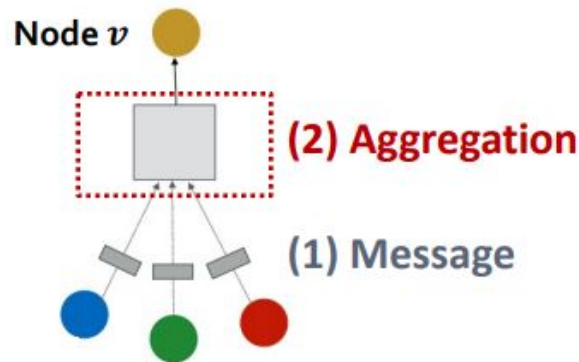
$$\mathbf{m}_u^{(l)} = \text{MSG}^{(l)} \left(\mathbf{h}_u^{(l-1)} \right), u \in \{N(v) \cup v\}$$

- **(2) Aggregation:** aggregate messages from neighbors

$$\mathbf{h}_v^{(l)} = \text{AGG}^{(l)} \left(\left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\}, \mathbf{m}_v^{(l)} \right)$$

- **Nonlinearity (activation)**

- Often written as $\sigma(\cdot)$: $\text{ReLU}(\cdot)$, $\text{Sigmoid}(\cdot)$, ...
- Can be added to message or aggregation



GNN layer
= Message (transformation) + aggregation

A single GNN layer

Summary:

- **(1) Message:** each node computes its own message

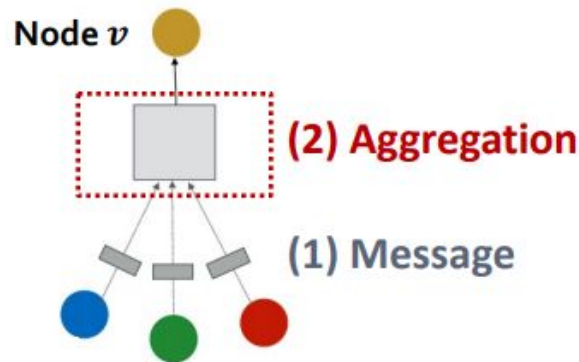
$$\mathbf{m}_u^{(l)} = \text{MSG}^{(l)}(\mathbf{h}_u^{(l-1)}), u \in \{N(v) \cup v\}$$

- **(2) Aggregation:** aggregate messages from neighbors

$$\mathbf{h}_v^{(l)} = \text{AGG}^{(l)}(\{\mathbf{m}_u^{(l)}, u \in N(v)\}, \mathbf{m}_v^{(l)})$$

- **Nonlinearity (activation)**

- Often written as $\sigma(\cdot)$: $\text{ReLU}(\cdot)$, $\text{Sigmoid}(\cdot)$, ...
- Can be added to message or aggregation



GNN layers have different instantiations

- GCN, GraphSAGE, GAT, ...
- Each has its own design of **MSG()** & **ADD()**

Implemented GNN layers in PyG



PyG (PyTorch Geometric)

<https://www.pyg.org/>

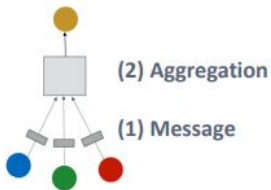
MessagePassing	Base class for creating message passing layers of the form
GCNConv	The graph convolutional operator from the "Semi-supervised Classification with Graph Convolutional Networks" paper
ChebConv	The chebyshev spectral graph convolutional operator from the "Convolutional Neural Networks on Graphs with Fast Localized Spectral Filtering" paper
SAGEConv	The GraphSAGE operator from the "Inductive Representation Learning on Large Graphs" paper
GraphConv	The graph neural network operator from the "Weisfeiler and Leman Go Neural: Higher-order Graph Neural Networks" paper
GravNetConv	The GravNet operator from the "Learning Representations of Irregular Particle-detector Geometry with Distance-weighted Graph Networks" paper, where the graph is dynamically constructed using nearest neighbors.
GatedGraphConv	The gated graph convolution operator from the "Gated Graph Sequence Neural Networks" paper
ResGatedGraphConv	The residual gated graph convolutional operator from the "Residual Gated Graph ConvNets" paper
GATConv	The graph attentional operator from the "Graph Attention Networks" paper
GATv2Conv	The GATv2 operator from the "How Attentive are Graph Attention Networks?" paper, which fixes the static attention problem of the standard <code>GATConv</code> layer: since the linear layers in the standard GAT are applied right after each other, the ranking of attended nodes is unconditioned on the query node.
TransformerConv	The graph transformer operator from the "Masked Label Prediction: Unified Message Passing Model for Semi-Supervised Classification" paper
AGNNConv	The graph attentional propagation layer from the "Attention-based Graph Neural Network for Semi-Supervised Learning" paper
TAGConv	The topology adaptive graph convolutional networks operator from the "Topology Adaptive Graph Convolutional Networks" paper
GINConv	The graph isomorphism operator from the "How Powerful are Graph Neural Networks?" paper
GINEConv	The modified <code>GINConv</code> operator from the "Strategies for Pre-training Graph Neural Networks" paper

Classical GNN layers: GCN

Graph convolutional networks (GCN)

$$\mathbf{h}_v^{(l)} = \sigma \left(\underbrace{\sum_{u \in N(v)} \mathbf{w}^{(l)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}}_{\text{Aggregation}} \right)$$

Message



■ Message:

- Each Neighbor: $\mathbf{m}_u^{(l)} = \frac{1}{|N(v)|} \mathbf{w}^{(l)} \mathbf{h}_u^{(l-1)}$

Normalized by node degree

(In the GCN paper they use a slightly different normalization)

■ Aggregation:

- Sum over messages from neighbors, then apply activation
- $\mathbf{h}_v^{(l)} = \sigma \left(\text{Sum} \left(\left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\} \right) \right)$

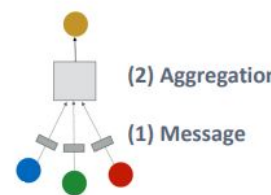
In GCN graph is assumed to have self-edges that are included in the summation.

Classical GNN layers: GCN

Graph convolutional networks (GCN)

$$\mathbf{h}_v^{(l)} = \sigma \left(\underbrace{\sum_{u \in N(v)} \mathbf{w}^{(l)} \frac{\mathbf{h}_u^{(l-1)}}{|N(v)|}}_{\text{Aggregation}} \right)$$

Message



All neighbors have the same weight in the average.
What if not all nodes are equally important?

■ Message:

- Each Neighbor: $\mathbf{m}_u^{(l)} = \frac{1}{|N(v)|} \mathbf{w}^{(l)} \mathbf{h}_u^{(l-1)}$

Normalized by node degree

(In the GCN paper they use a slightly different normalization)

■ Aggregation:

- Sum** over messages from neighbors, then apply activation
- $\mathbf{h}_v^{(l)} = \sigma \left(\text{Sum} \left(\left\{ \mathbf{m}_u^{(l)}, u \in N(v) \right\} \right) \right)$

In GCN graph is assumed to have self-edges that are included in the summation.

Classical GNN layers: GAT

Graph attention networks (GAT)

$$\mathbf{h}_v^{(l)} = \sigma(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

Attention weights

- **Solution: Attention mechanism**
 - **Goal:** Learn weights α_{vu} from data, instead of specifying manually (e.g., $1/N(v)$)
 - Used as a drop-in layer to aggregate embeddings in a neural network (not only in GNN)

Attention mechanism

(1) Apply small neural network (e.g., a single layer) a to compute the **attention coefficients** e_{vu} across pairs of nodes u, v based on their messages:

$$e_{vu} = a(\mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)}, \mathbf{W}^{(l)} \mathbf{h}_v^{(l-1)})$$

- e_{vu} indicates the importance of u 's message to node v

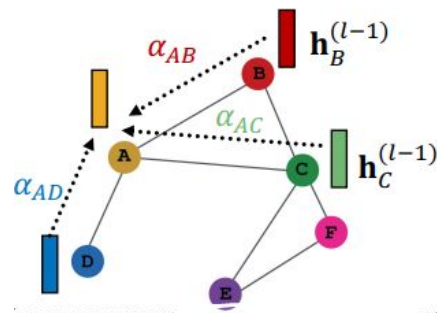
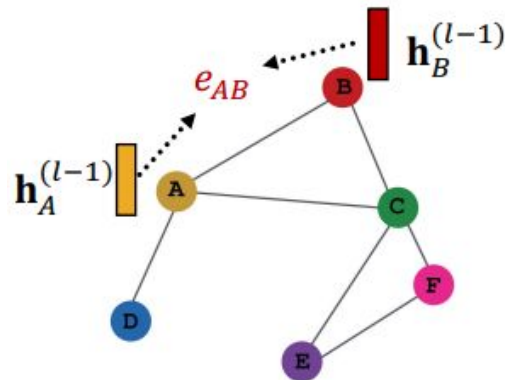
(2) **Normalize** e_{vu} into the final attention weight α_{vu}

- Use the **softmax** function, so that $\sum_{u \in N(v)} \alpha_{vu} = 1$:

$$\alpha_{vu} = \frac{\exp(e_{vu})}{\sum_{k \in N(v)} \exp(e_{vk})}$$

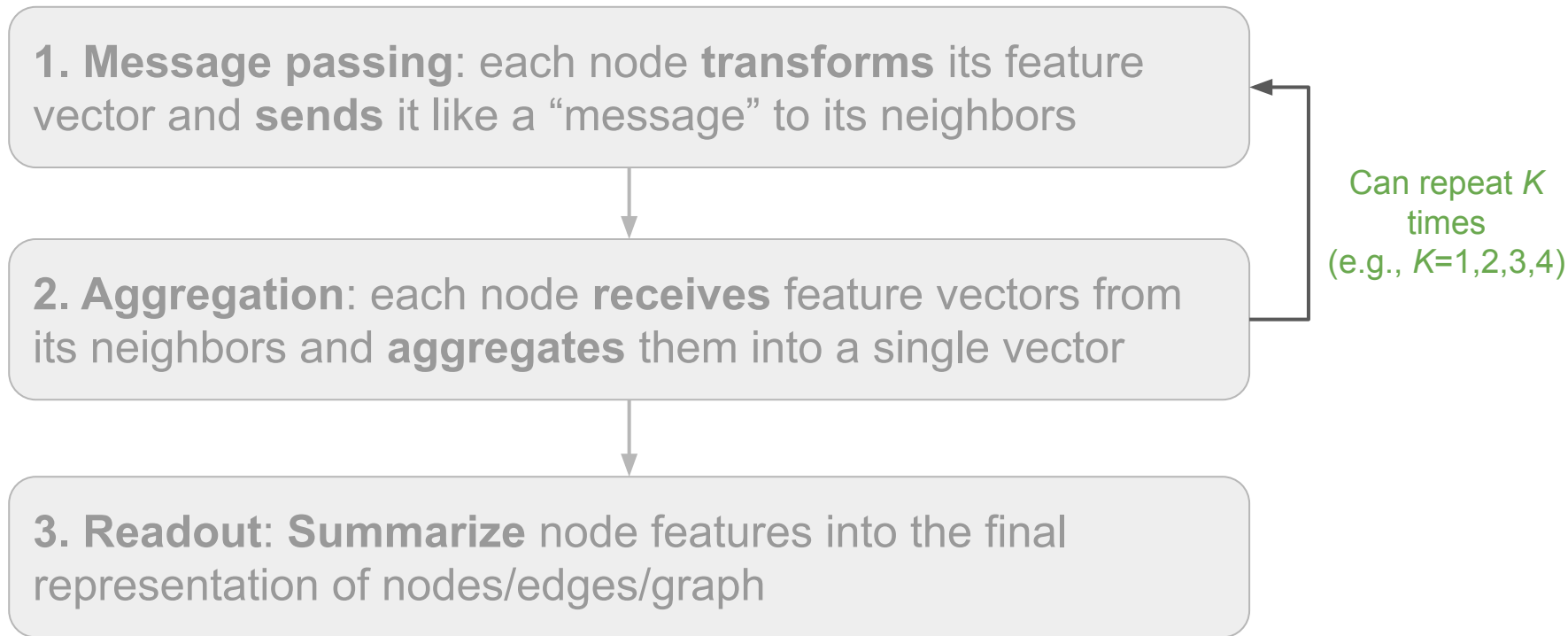
(3) **Weighted sum** based on the final attention weight α_{vu}

$$\mathbf{h}_v^{(l)} = \sigma(\sum_{u \in N(v)} \alpha_{vu} \mathbf{W}^{(l)} \mathbf{h}_u^{(l-1)})$$

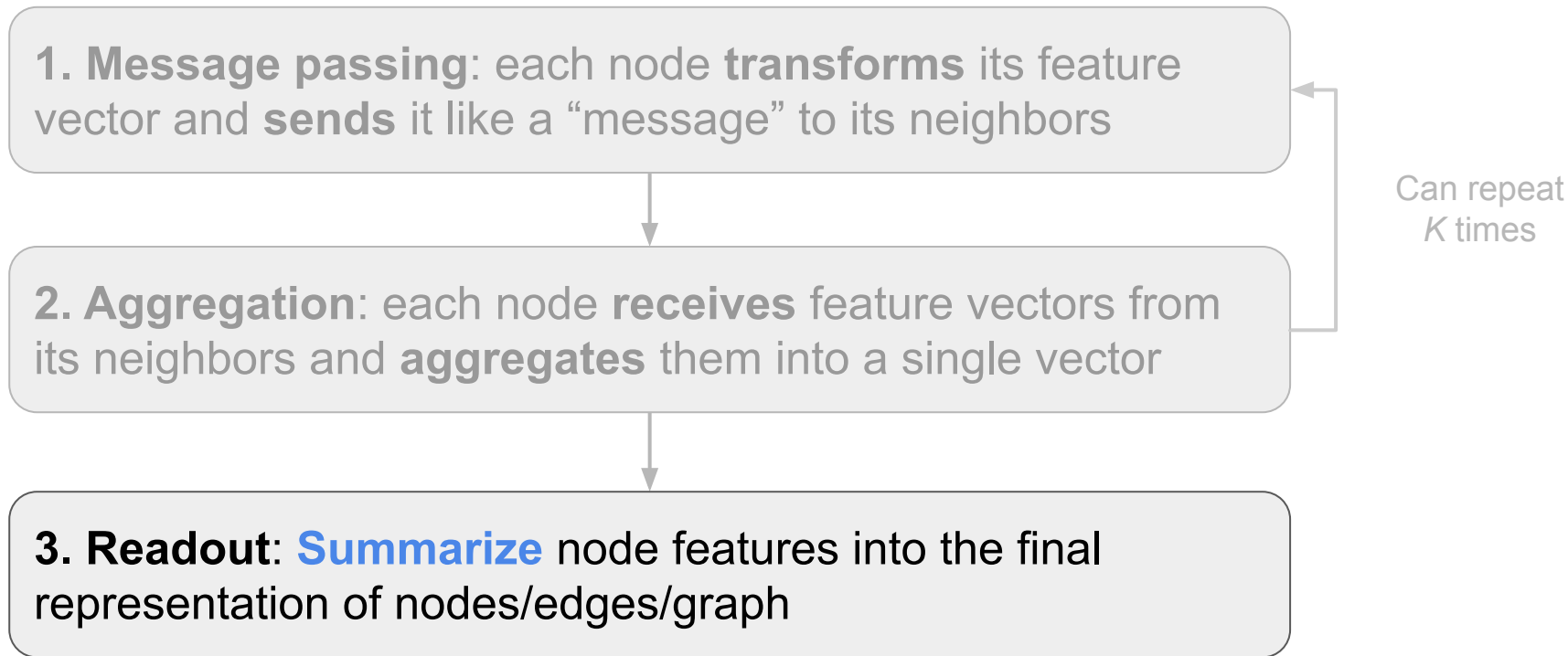


Weighted sum using $\alpha_{AB}, \alpha_{AC}, \alpha_{AD}$:
 $\mathbf{h}_A^{(l)} = \sigma(\alpha_{AB} \mathbf{W}^{(l)} \mathbf{h}_B^{(l-1)} + \alpha_{AC} \mathbf{W}^{(l)} \mathbf{h}_C^{(l-1)} + \alpha_{AD} \mathbf{W}^{(l)} \mathbf{h}_D^{(l-1)})$

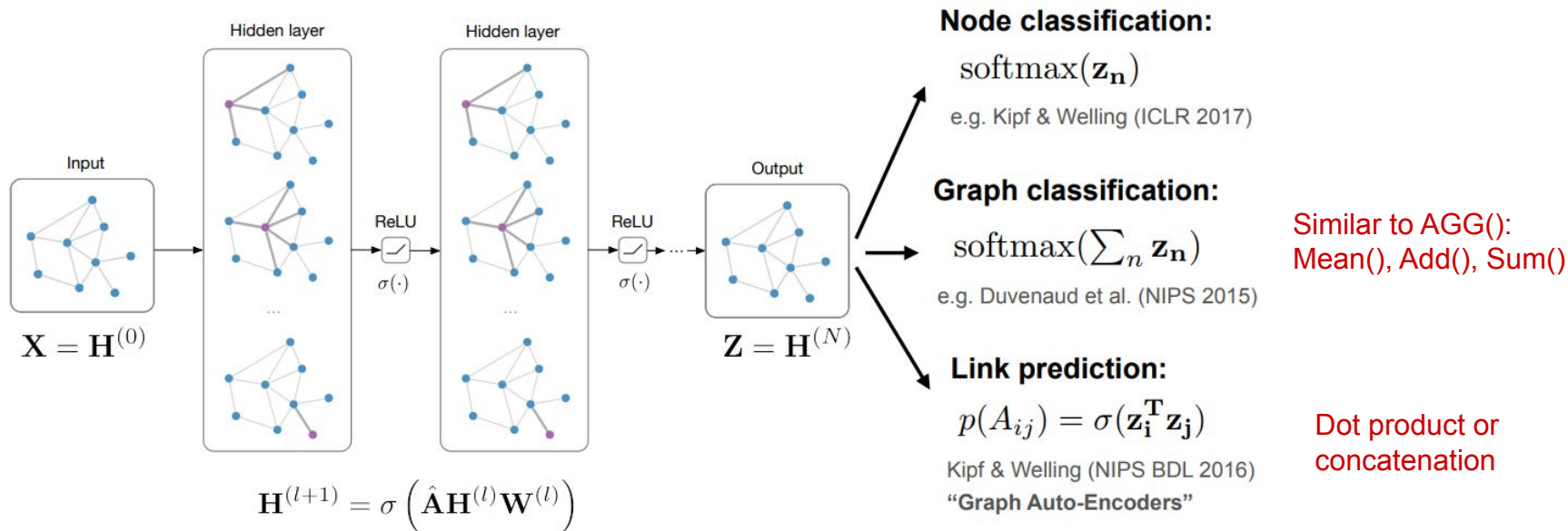
General GNN framework



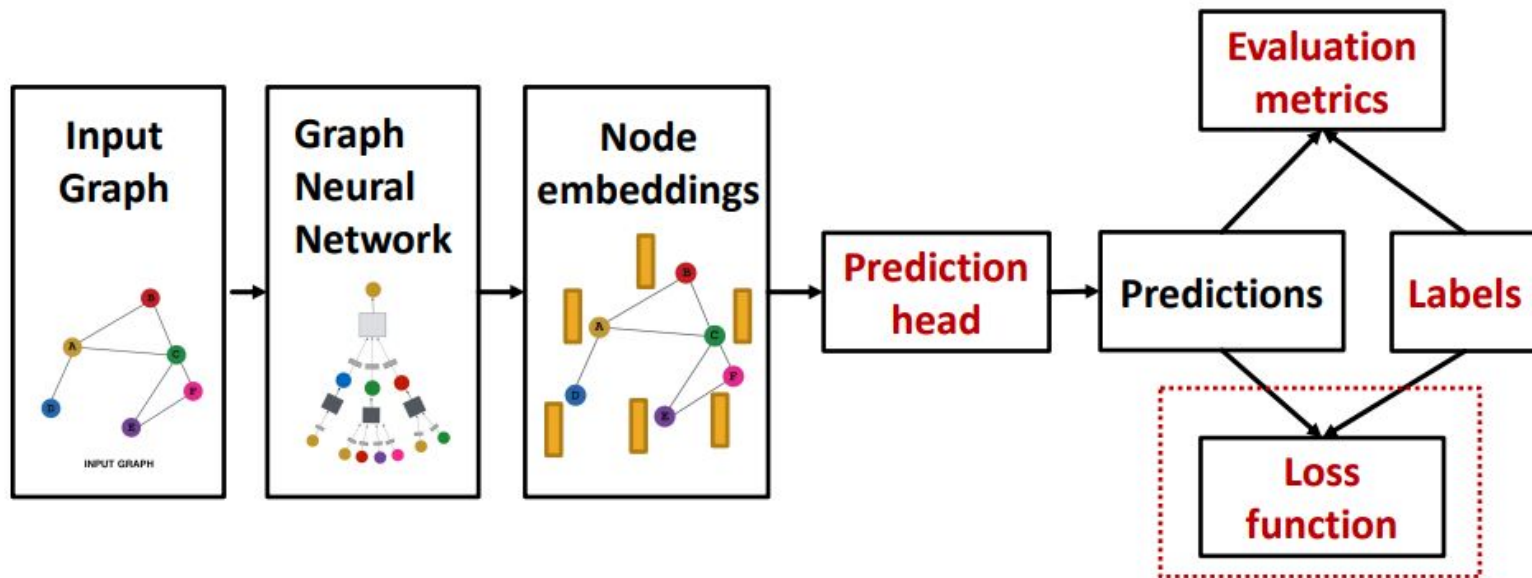
General GNN framework



Readout operation depends on the task



Training GNN models



Classification & Regression

- **Classification:** labels $y^{(i)}$ with **discrete** value
 - E.g., Node classification: which category does a node belong to
- **Regression:** labels $y^{(i)}$ with **continuous** value
 - E.g., predict the drug likeness of a molecular graph
- GNNs can be applied to both settings
- **Differences:** **loss function** & **evaluation metrics**

Classification loss

- **Cross entropy (CE)** is a common loss function in classification

- **K -way prediction** for i -th data point:

$$\text{CE}(\underbrace{\mathbf{y}^{(i)}}_{\text{Label}}, \underbrace{\hat{\mathbf{y}}^{(i)}}_{\text{Prediction}}) = - \sum_{j=1}^K \mathbf{y}_j^{(i)} \log(\hat{\mathbf{y}}_j^{(i)})$$

i -th data point
 j -th class

where:

E.g.

0	0	1	0	0
---	---	---	---	---

$\mathbf{y}^{(i)} \in \mathbb{R}^K$ = one-hot label encoding

$\hat{\mathbf{y}}^{(i)} \in \mathbb{R}^K$ = prediction after Softmax(\cdot)

E.g.

0.1	0.3	0.4	0.1	0.1
-----	-----	-----	-----	-----

- **Total loss over all N training examples**

$$\text{Loss} = \sum_{i=1}^N \text{CE}(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

Regression loss

- **Mean squared loss (MSE)** is a common loss function in regression

- **K -way regression** for data point (i):

$K=1$ in most of the applications

$$\text{MSE}(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)}) = \sum_{j=1}^K (\mathbf{y}_j^{(i)} - \hat{\mathbf{y}}_j^{(i)})^2$$

i -th data point
 j -th target

where:

E.g.

1.4	2.3	1.0	0.5	0.6
-----	-----	-----	-----	-----

$\mathbf{y}^{(i)} \in \mathbb{R}^k$ = Real valued vector of targets
 $\hat{\mathbf{y}}^{(i)} \in \mathbb{R}^k$ = Real valued vector of predictions

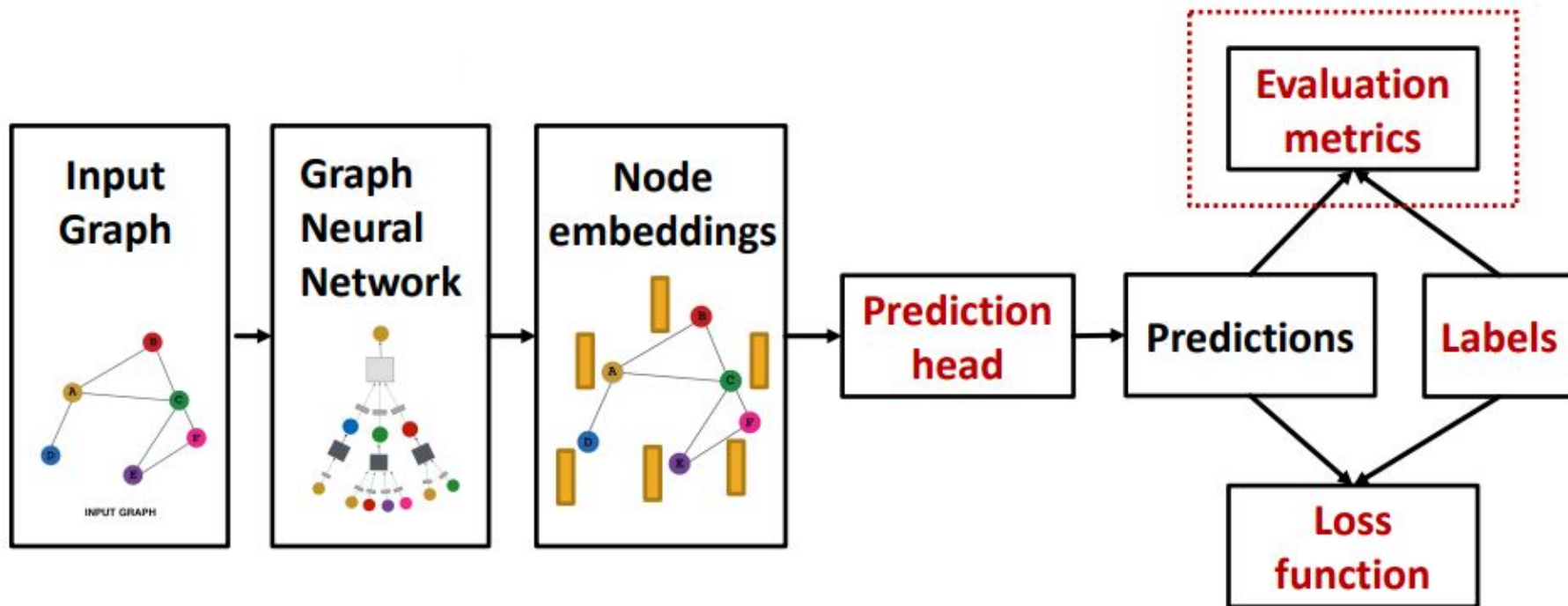
E.g.

0.9	2.8	2.0	0.3	0.8
-----	-----	-----	-----	-----

- **Total loss over all N training examples**

$$\text{Loss} = \sum_{i=1}^N \text{MSE}(\mathbf{y}^{(i)}, \hat{\mathbf{y}}^{(i)})$$

GNN evaluation



Regression metrics

- Root mean square error (RMSE)

$$\sqrt{\sum_{i=1}^N \frac{(\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)})^2}{N}}$$

- Mean absolute error (MAE)

$$\frac{\sum_{i=1}^N |\mathbf{y}^{(i)} - \hat{\mathbf{y}}^{(i)}|}{N}$$

$\hat{\mathbf{y}}^{(i)}$: model prediction

$\mathbf{y}^{(i)}$: true label

Classification metrics

- **Binary classification:**

- **Accuracy:**

$$\frac{TP + TN}{TP + TN + FP + FN} = \frac{TP + TN}{|\text{Dataset}|}$$

- **Precision (P):**

$$\frac{TP}{TP + FP}$$

- **Recall (R):**

$$\frac{TP}{TP + FN}$$

- **F1-Score:**

$$\frac{2P * R}{P + R}$$

Confusion matrix

	Actually Positive (1)	Actually Negative (0)
Predicted Positive (1)	True Positives (TPs)	False Positives (FPs)
Predicted Negative (0)	False Negatives (FNs)	True Negatives (TNs)

- **Multi-class classification:**

Accuracy

$$\frac{1[\text{argmax}(\hat{\mathbf{y}}^{(i)}) = \mathbf{y}^{(i)}]}{N}$$

Classification metric: ROC

- **ROC Curve:** Captures the tradeoff in TPR and FPR as the classification threshold is varied for a **binary** classifier.

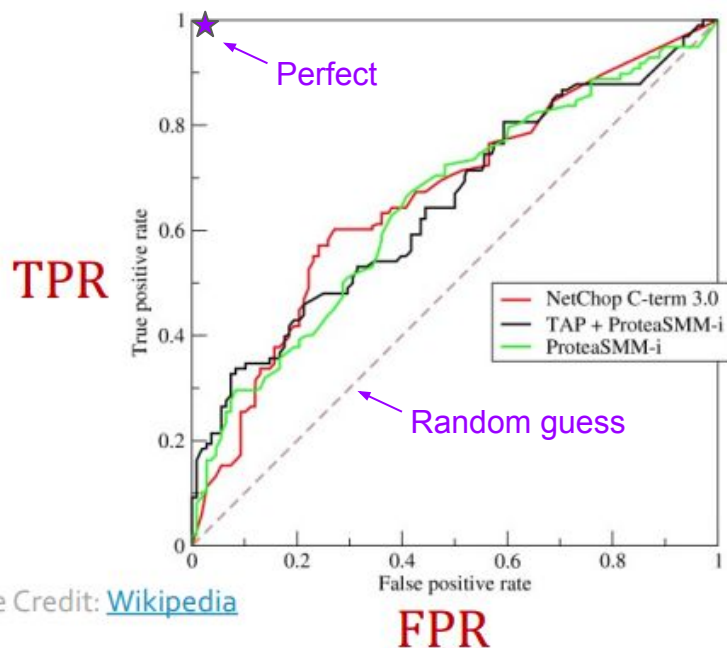


Image Credit: [Wikipedia](#)

$$\text{TPR} = \text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}$$

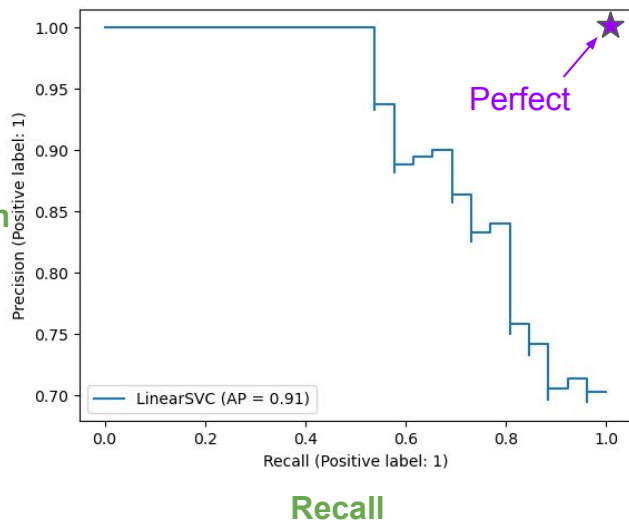
$$\text{FPR} = \frac{\text{FP}}{\text{FP} + \text{TN}}$$

- **Metric: AUROC** (Area under the ROC Curve)
 - 1.0: perfect prediction
 - 0.0: worst (random guess)

What if you get a classifier with AUROC = 0.1?

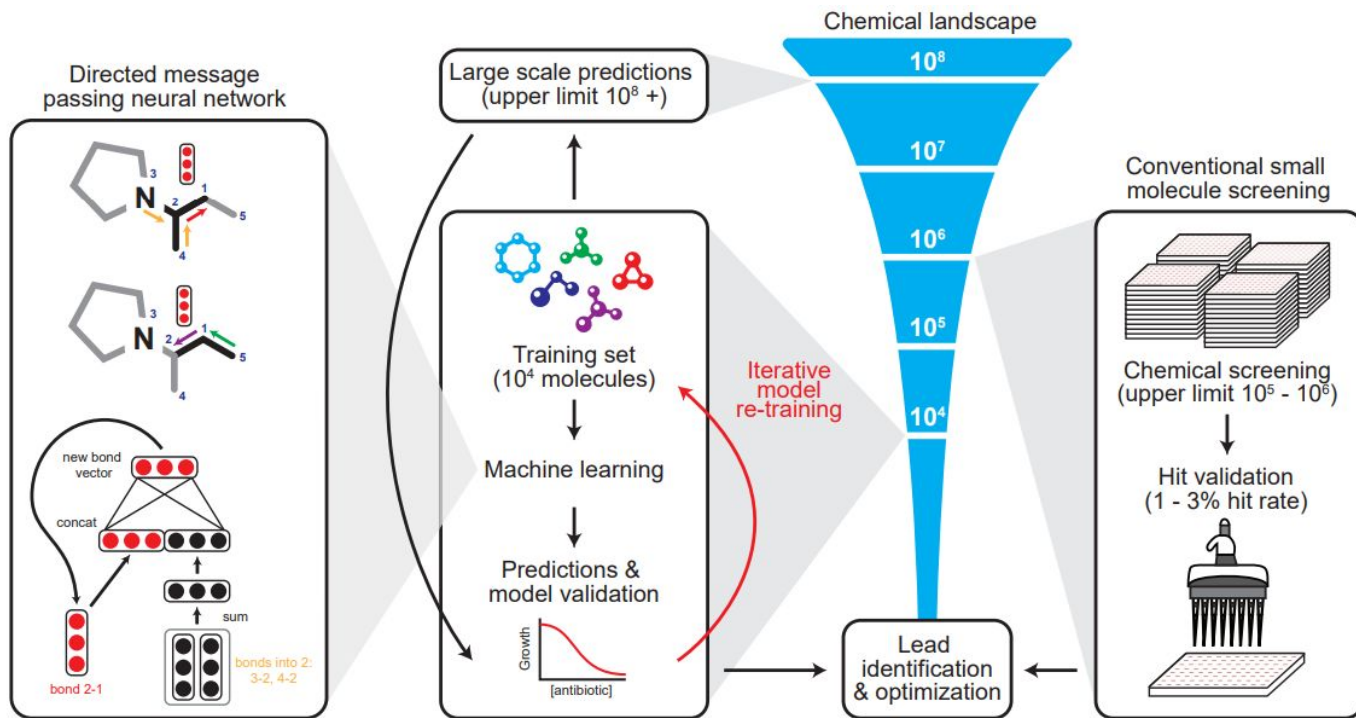
Classification metric: Precision-Recall curve

- **PR Curve:** Captures the tradeoff in **Precision** and **Recall** as the classification threshold is varied for a **binary** classifier.

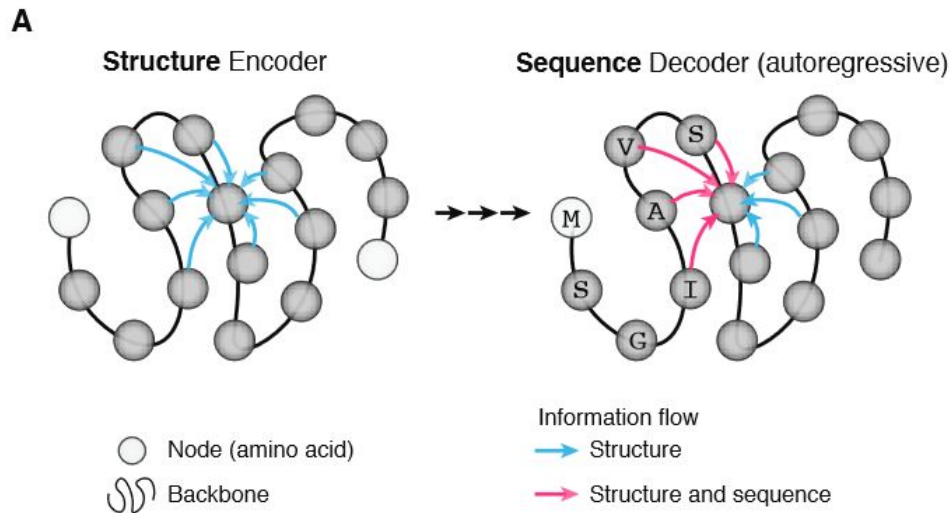


- **Metric: AUPR** (Area under the PR curve)
 - 1.0: perfect prediction
 - 0.0: worst

Applications: antibiotic discovery



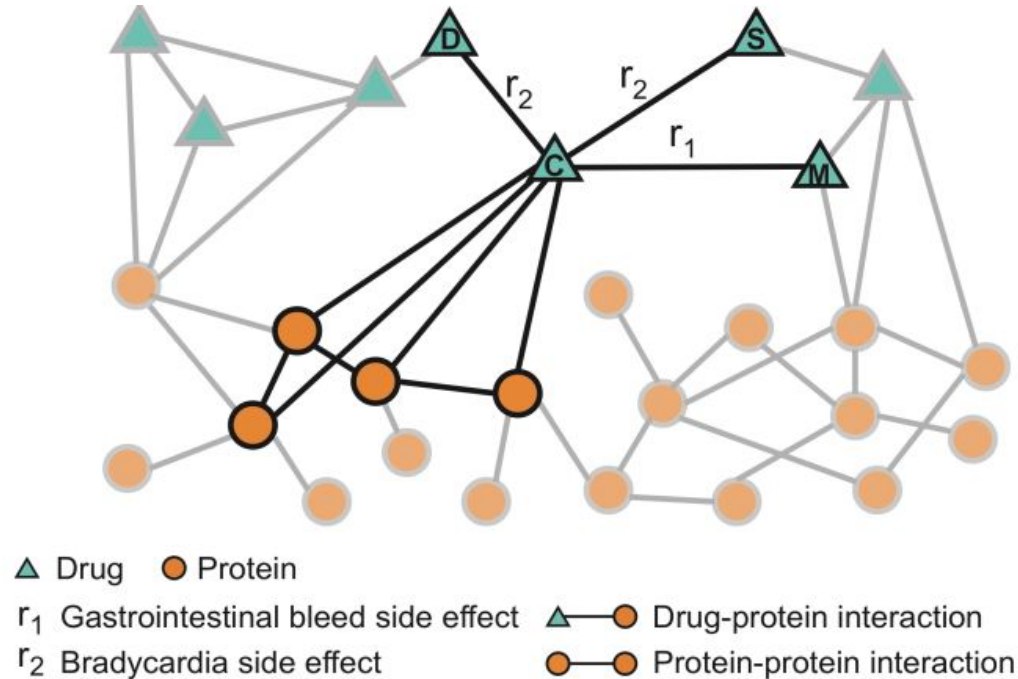
Applications: protein sequence design



Ingraham et al. "Generative models for graph-based protein design", NeurIPS, 2019

Jin et al. "Iterative Refinement Graph Neural Network for Antibody Sequence-Structure Co-Design", ICLR 2022

Applications: Polypharmacy effect prediction



Summary of today

- Graph neural network (GNN)
 - Generalize convolution to graphs
 - Invariance and equivariance
- GNN framework
 - Message
 - Aggregation
 - Readout