

Graph Neural Networks



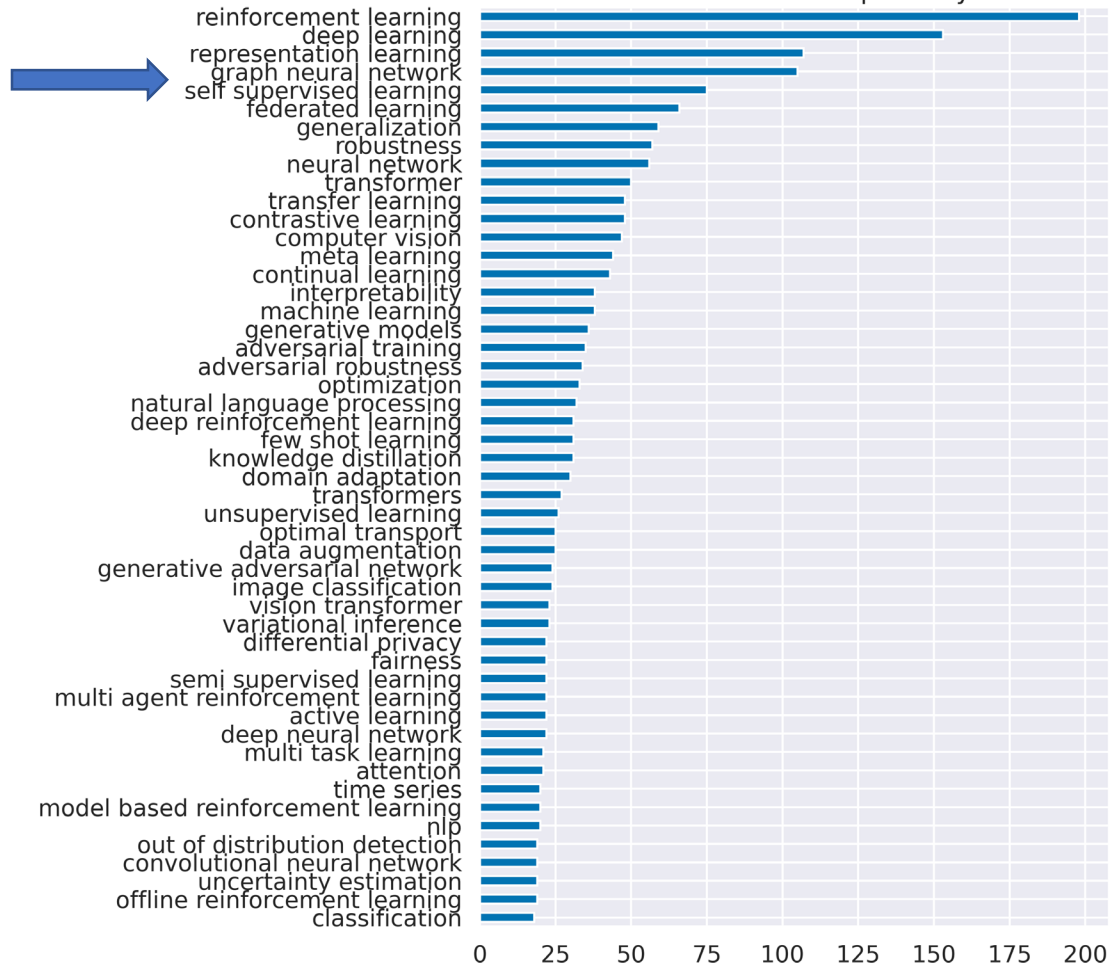
上海交通大学
约翰·霍普克罗夫特
计算机科学中心

John Hopcroft Center for Computer Science

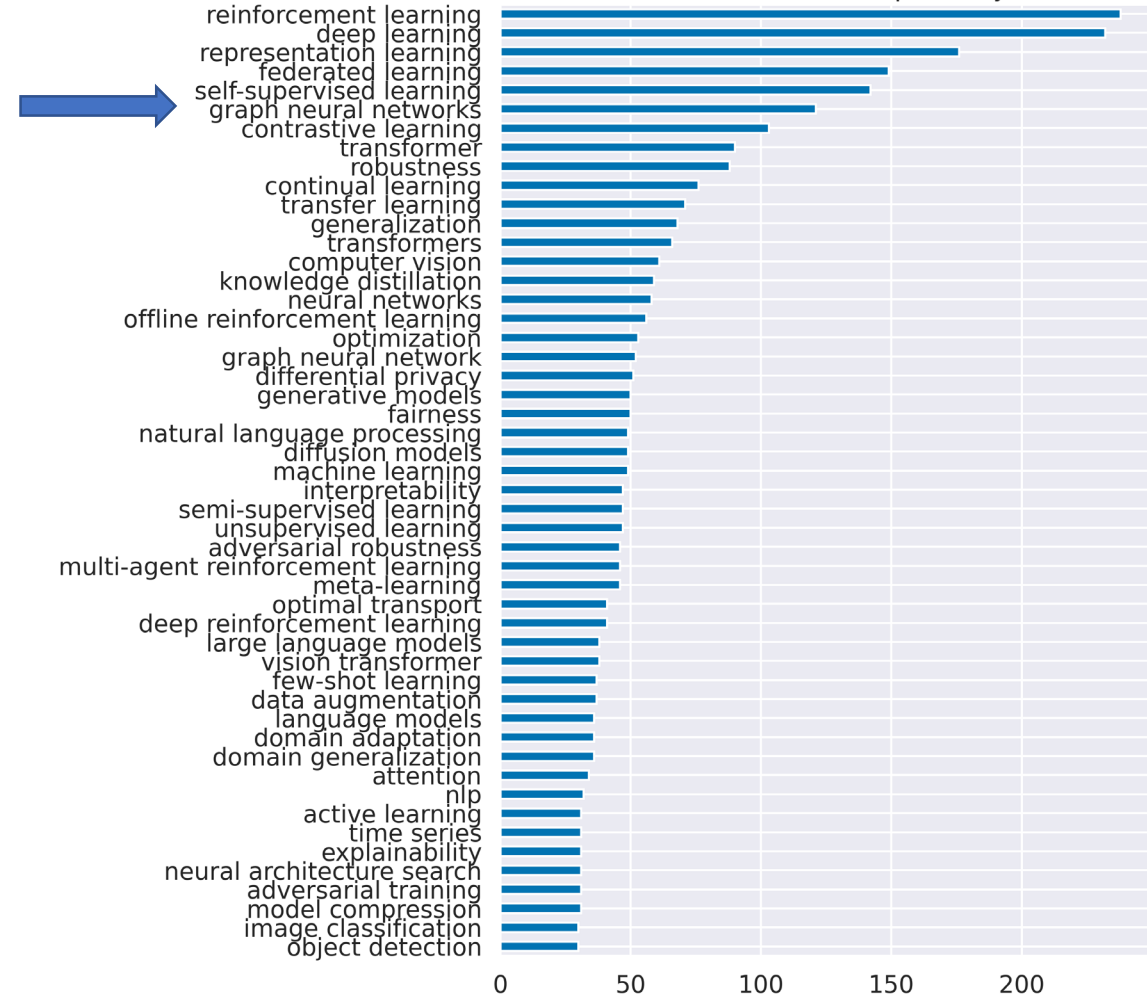


GNN is a hot topic in recent years

ICLR 2022 Submission Top 50 Keywords

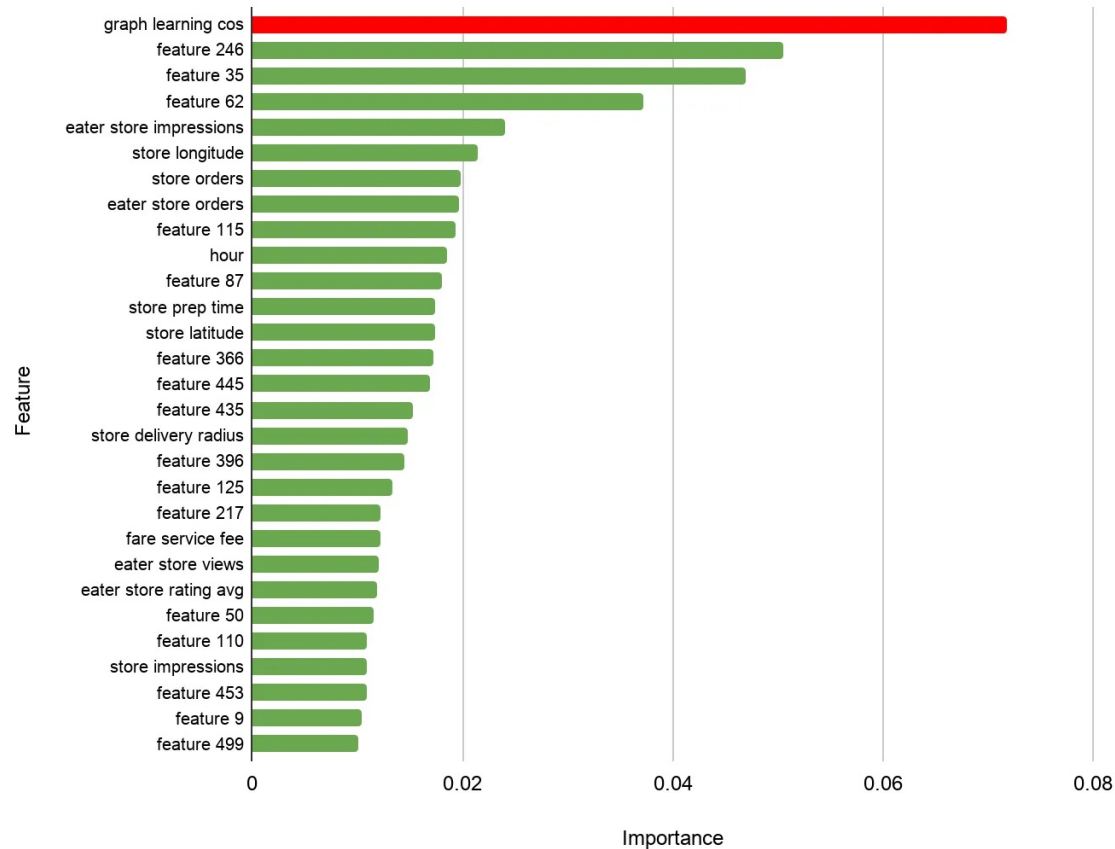


ICLR 2023 Submission Top 50 Keywords

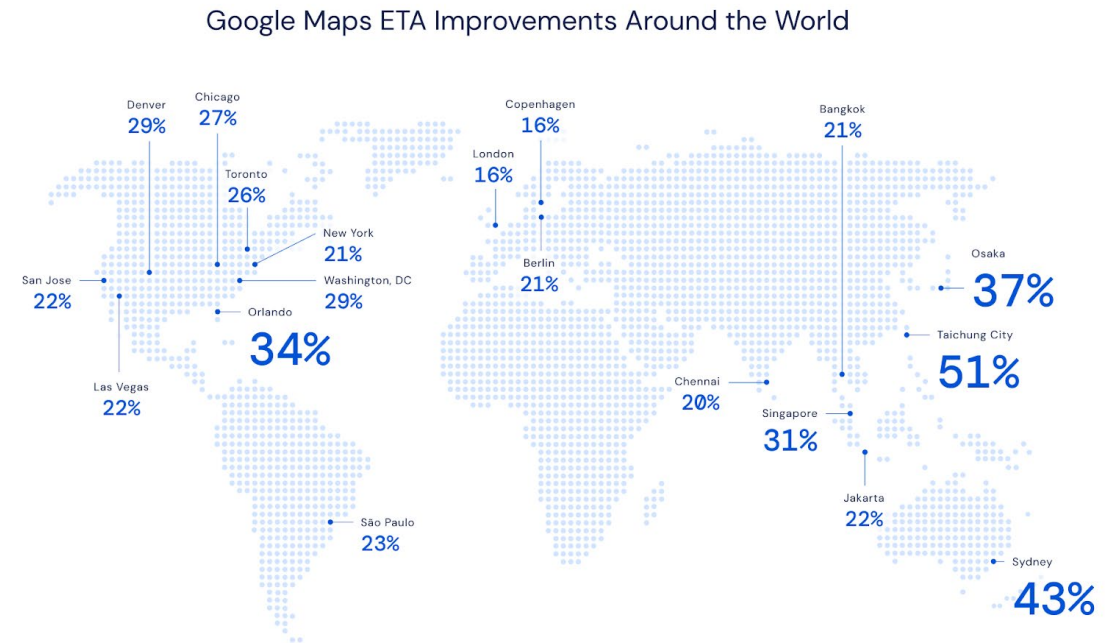


Applications

Uber Eats recommendation system 20% ↑

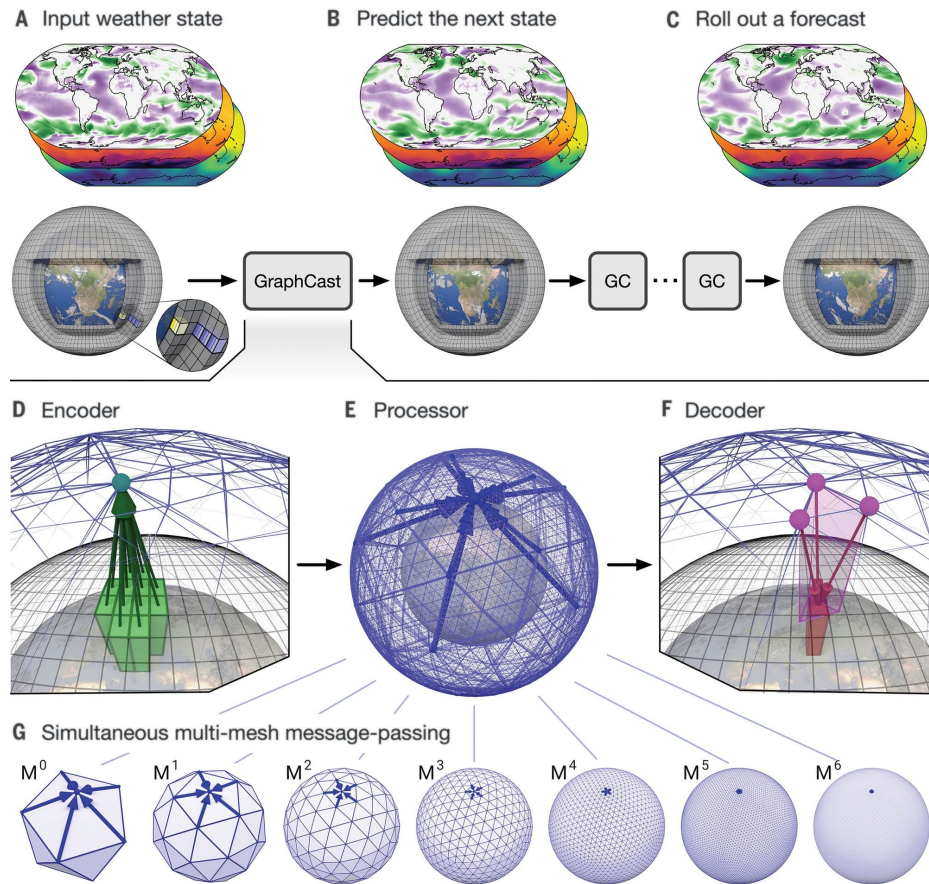


Google Maps ETA accuracy 50% ↑

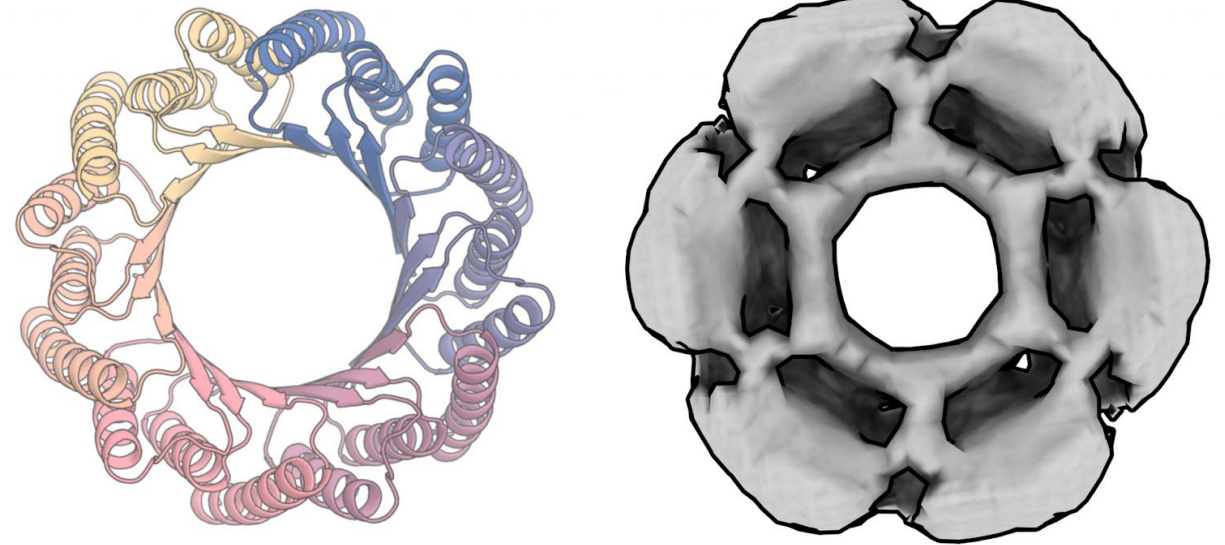


Applications

DeepMind GraphCast Climate Prediction



Baker Lab Protein Design 5x improvement



<https://www.science.org/doi/10.1126/science.adi2336?ref=assemblyai.com>

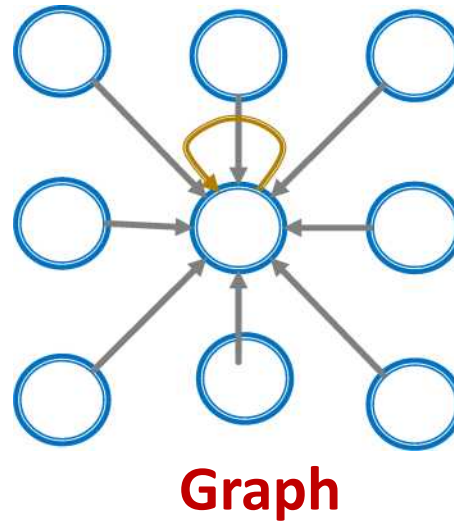
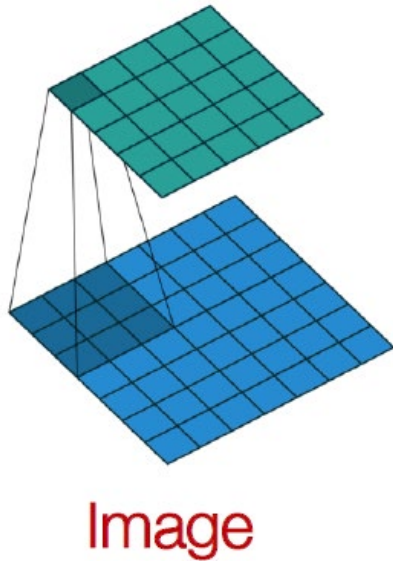
https://www.bakerlab.org/wp-content/uploads/2022/11/Diffusion_preprint_12012022.pdf?ref=assemblyai.com

Setup

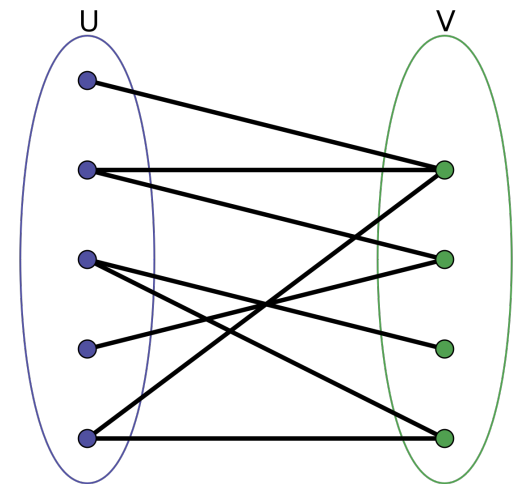
- **Graph G :**
- V is the **vertex set**
 - v : a **node** in V ;
 - $N(v)$: the set of **neighbors** of v
- A is the **adjacency matrix** (assume binary)
- $X \in \mathbb{R}^{m \times |V|}$ is a matrix of **node features**
 - Social networks: user profile
 - Biological networks: gene functional information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

Homophily: birds of a feather flock together

- **Idea:** transform information at the **neighbors** and combine it:
 - Transform “messages” h_i from neighbors: $W_i h_i$
 - Add them up: $\sum_i W_i h_i$

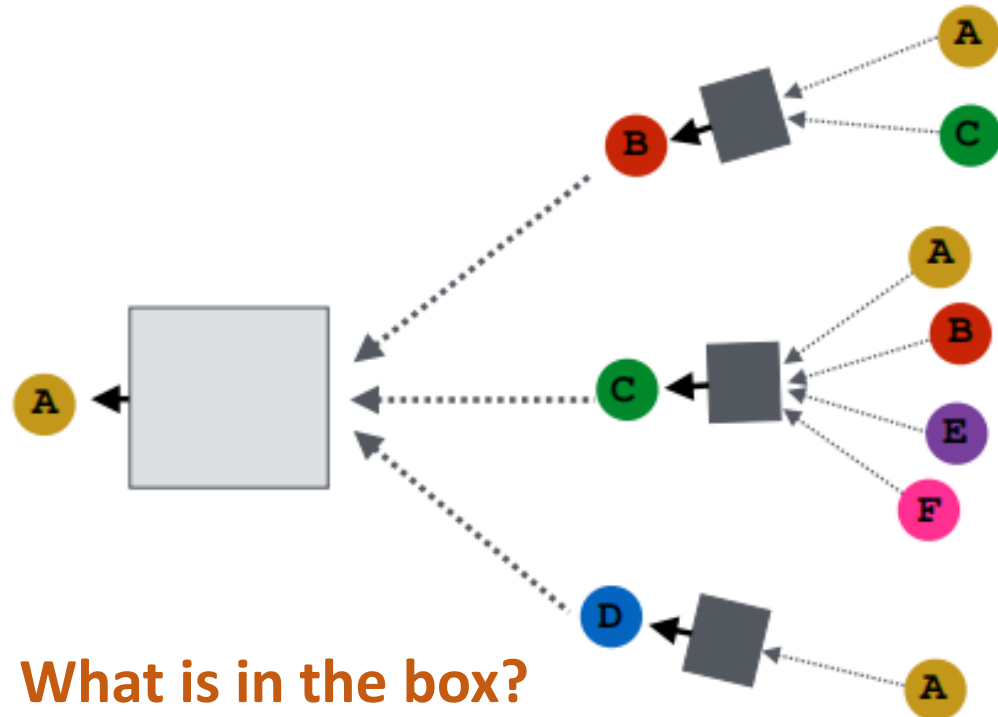
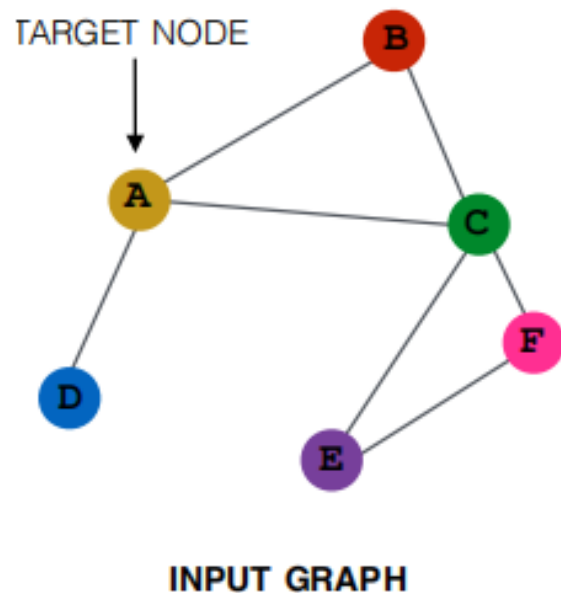


Example for the least homophily



Idea: Aggregate Neighbors

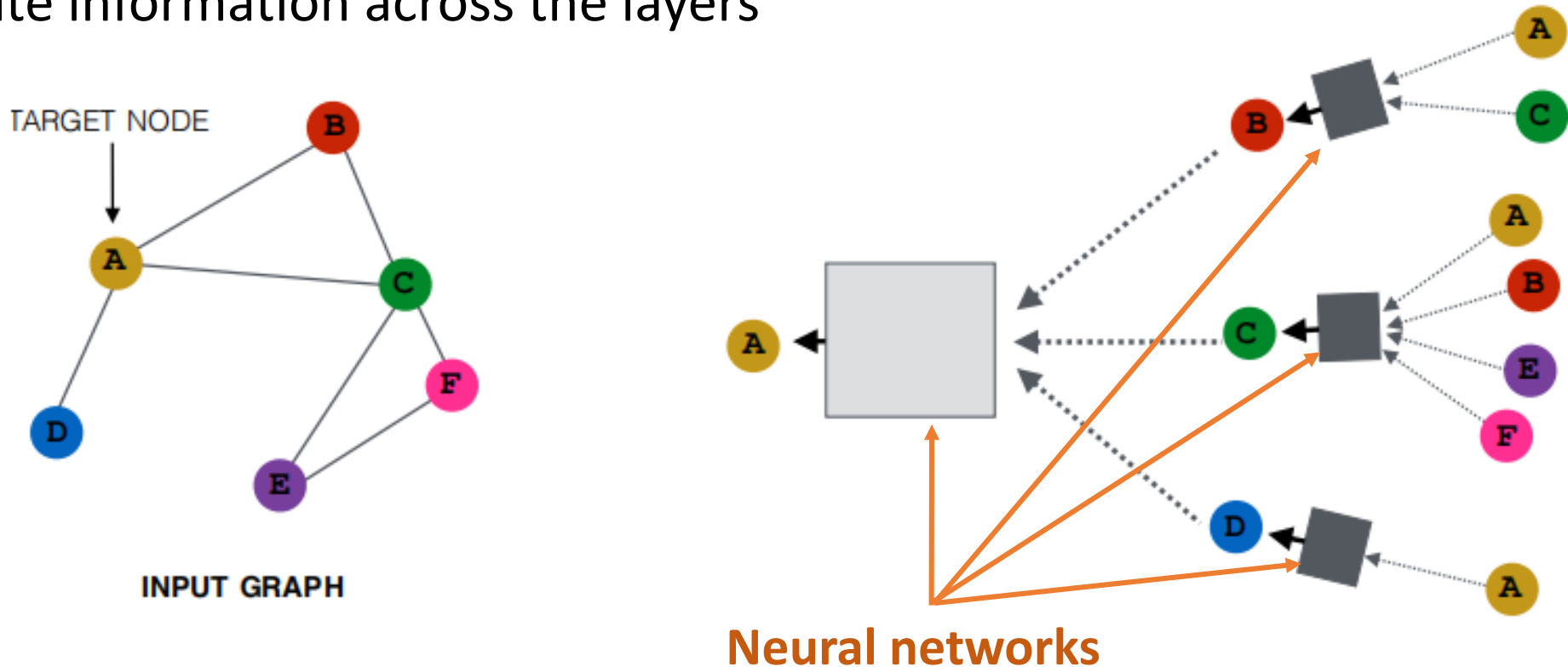
- **Key idea:** Generate node embeddings based on **local network neighborhoods**



Order insensitive functions: mean, max, some NNs

Graph Neural Networks

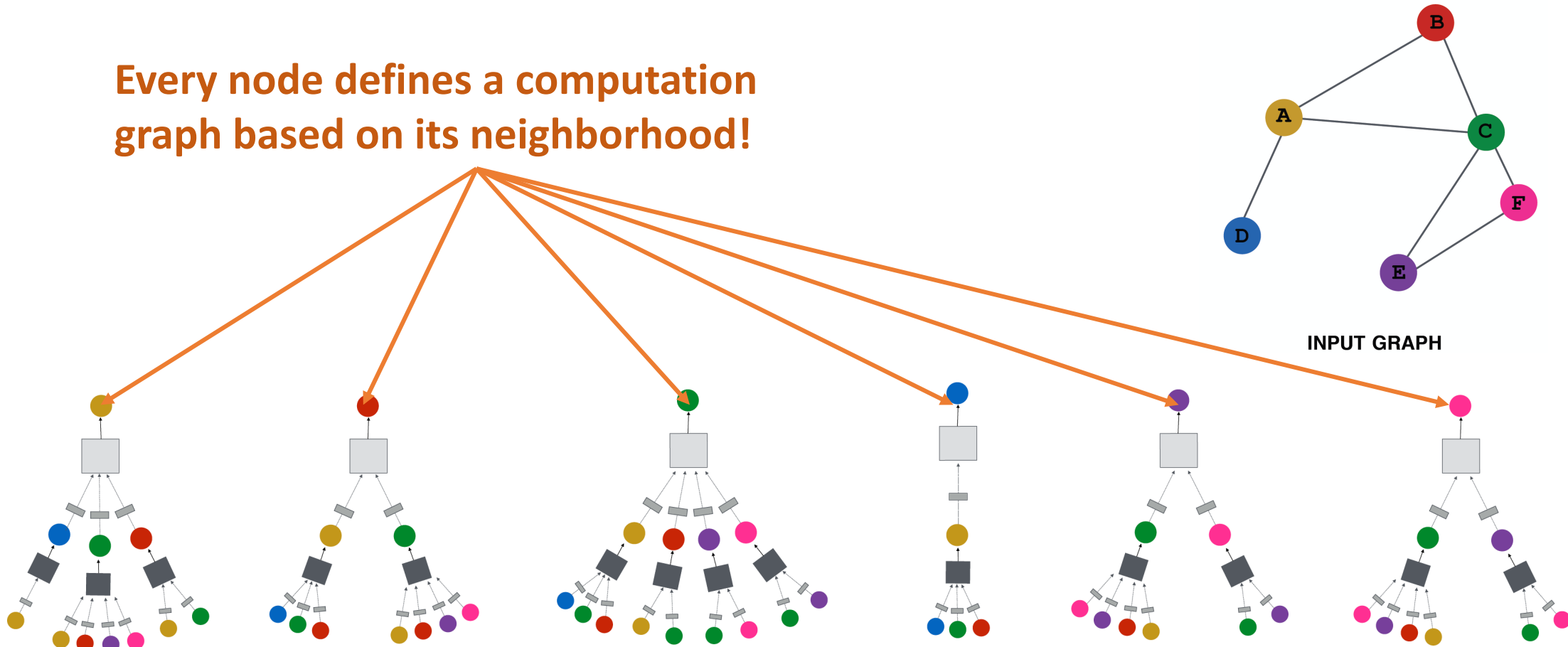
- **Intuition:** Nodes aggregate information from their neighbors using neural networks
- Key distinctions of different GNN are in how different approaches aggregate information across the layers



Graph Neural Networks

- **Intuition:** Network neighborhood defines a **computation graph**

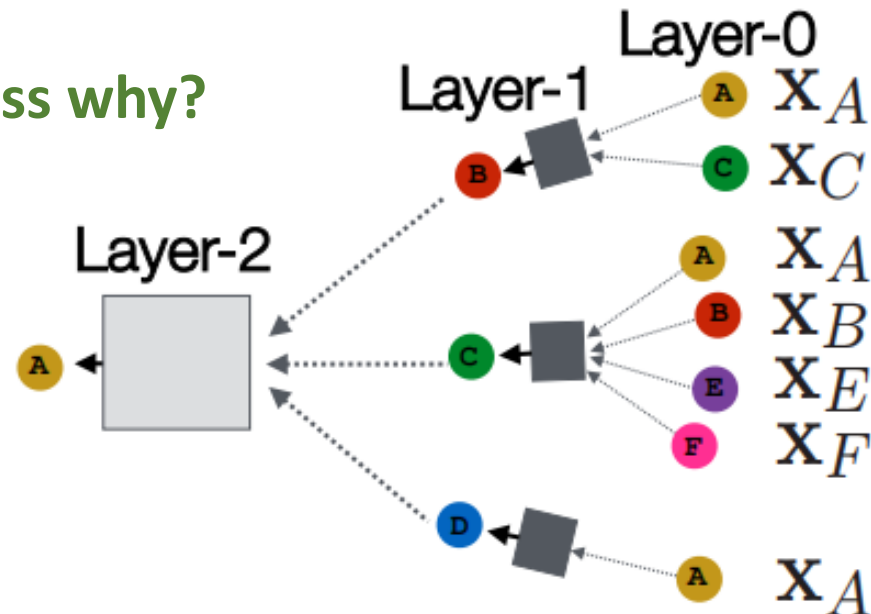
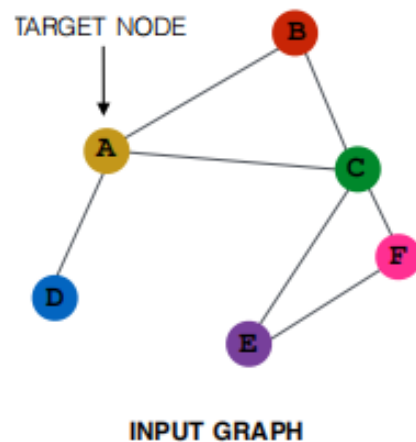
Every node defines a computation graph based on its neighborhood!



Model Layers

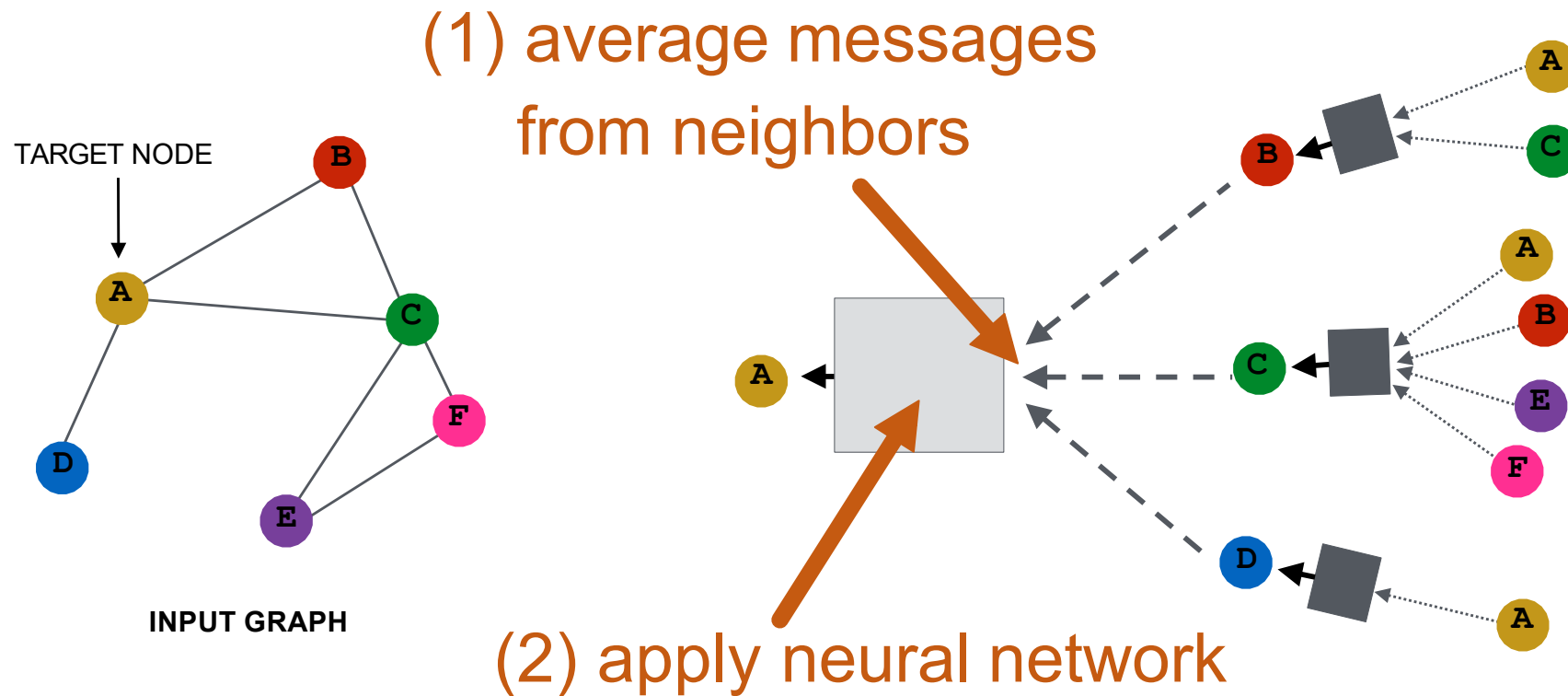
- Model can be **of arbitrary depth**:
 - Nodes have embeddings at each layer
 - Layer-0 embedding of node u is its input feature, x_u
 - Layer- k embedding gets information from nodes that are K hops away

But usually GNN can not be too deep. Guess why?

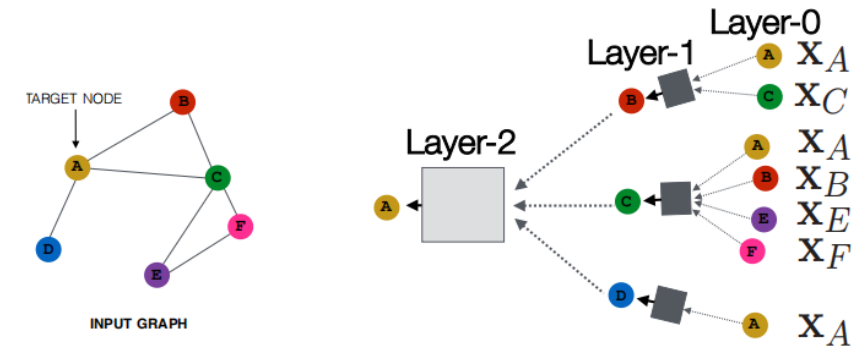


Basic Neighborhood Aggregation

- **Basic approach:** Average information from neighbors and apply a neural network



The Math: Deep Encoder



- **Basic approach:** Average neighbor messages and apply a neural network

$$h_v^0 = x_v$$

Initial 0-th layer embeddings are equal to node features

$$h_v^{(l+1)} = \sigma \left(W_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)} \right), \forall l \in \{0, \dots, L-1\}$$

embedding of v at layer l

$$z_v = h_v^{(L)}$$

Embedding after L-layers of neighborhood aggregation

Non-linearity (e.g., Sigmoid)

Average of neighbor's previous layer embeddings

Total number of layers

Model Parameters

Trainable weight matrices (i.e., what we learn)

$$h_v^{(0)} = x_v$$
$$h_v^{(l+1)} = \sigma \left(W_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)} \right), \forall l \in \{0, \dots, L-1\}$$

$$z_v = h_v^{(L)}$$

Final node embedding

We can feed these **embeddings into any loss function** and run **SGD** to train the **weight** parameters

- h_v^l : the hidden representation of node v at layer l
- W_l : weight matrix for neighborhood aggregation
- B_l : weight matrix for transforming hidden vector of self

Matrix Formulation

- Many aggregations can be performed efficiently by matrix operations

- Let $H^{(l)} = \begin{bmatrix} h_1^{(l)} & \dots & h_{|V|}^{(l)} \end{bmatrix}^T$

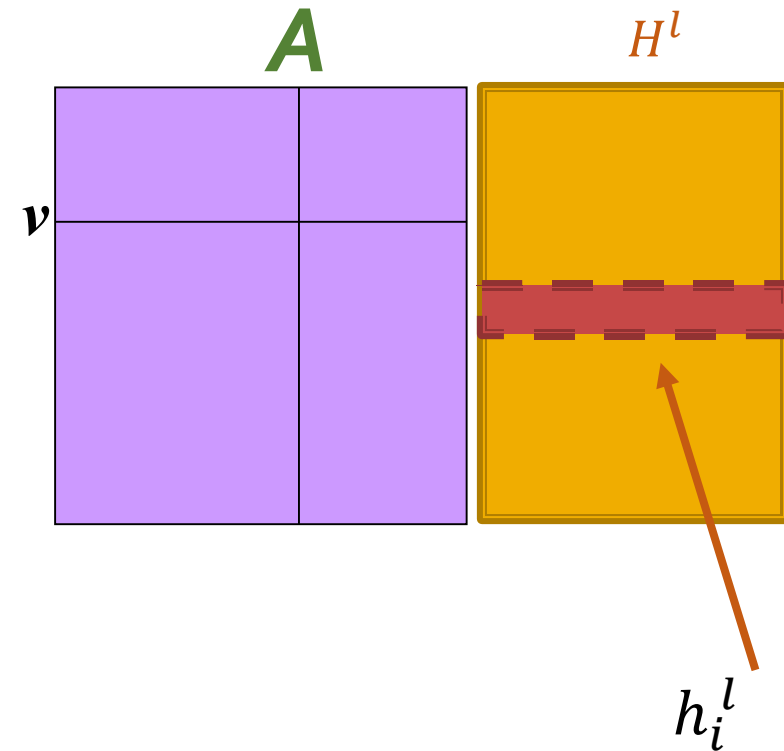
- Then: $\sum_{u \in N_v} h_u^{(l+1)} = A_{v,:} H^{(l)}$

- Let D be diagonal matrix where $D_{v,v} = \text{Deg}(v) = |N(v)|$

- The inverse of D : D^{-1} is also diagonal: $D_{v,v}^{-1} = 1/|N(v)|$

- Therefore,

$$\sum_{u \in N(v)} \frac{h_u^l}{|N(v)|} \longrightarrow H^{(l+1)} = D^{-1} A H^{(l)}$$



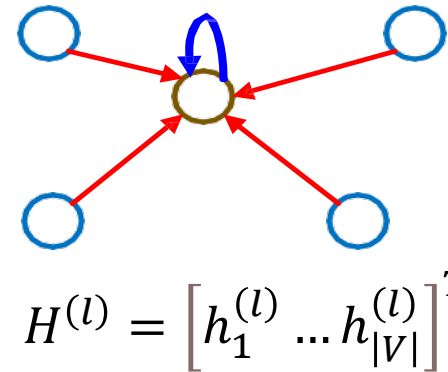
Matrix Formulation

- Re-writing update function in matrix form:

- $H^{(l+1)} = \sigma(\tilde{A}H^{(l)}W_l + H^{(l)}B_l)$ 

where $\tilde{A} = D^{-1}A$

- **Red:** neighborhood aggregation
- **Green:** self transformation



How to train a GNN

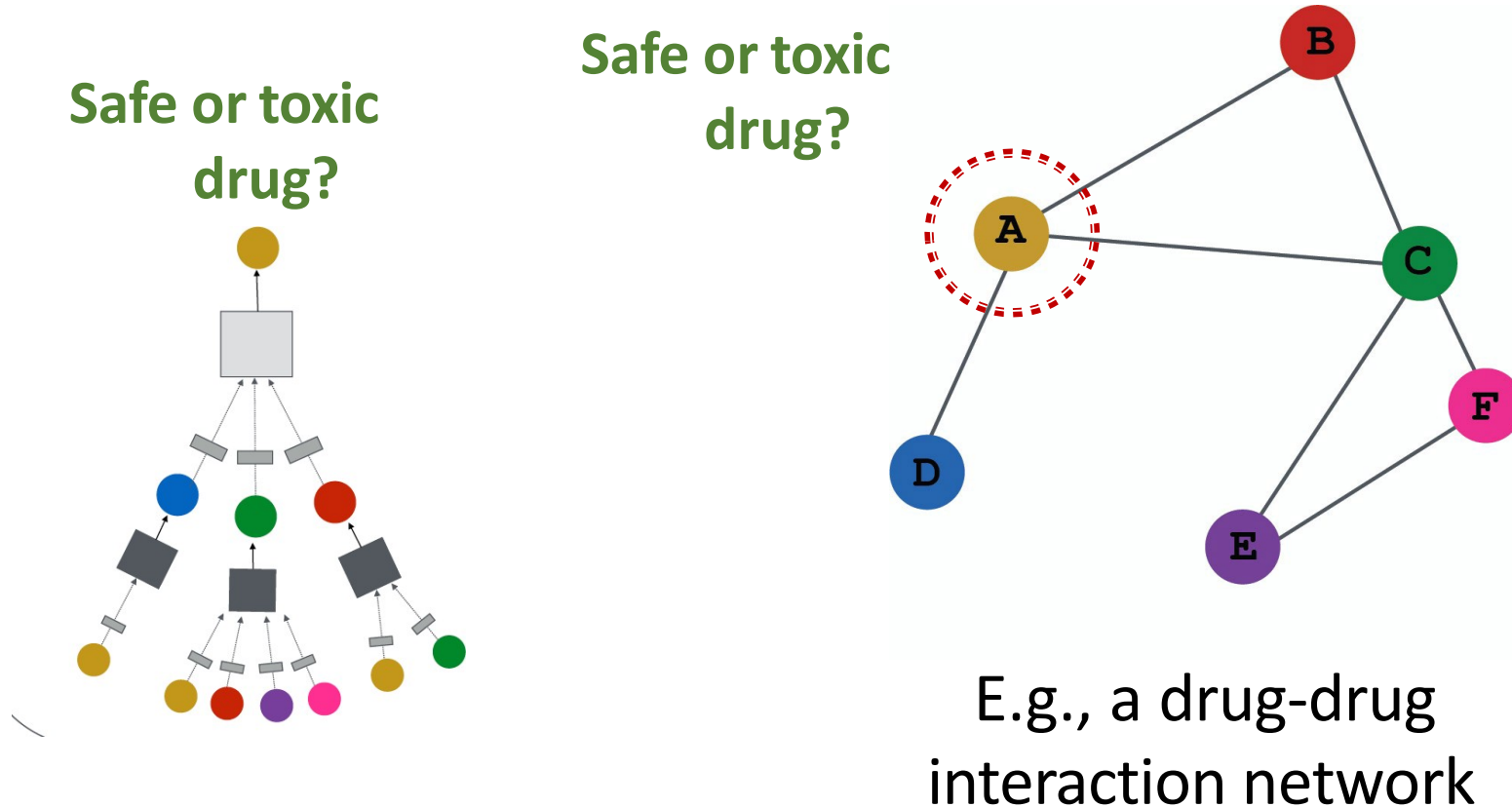
- Node embedding \mathbf{z}_v is a function of input graph
- **Supervised setting:** we want to minimize the loss \mathcal{L} :

$$\min_{\Theta} \mathcal{L}(\mathbf{y}, f(\mathbf{z}_v))$$

- \mathbf{y} : node label
 - \mathcal{L} could be L2 if \mathbf{y} is real number, or cross entropy if \mathbf{y} is categorical
-
- **Unsupervised setting:**
 - No node label available
 - **Use the graph structure as the supervision!**

Supervised Training

- **Directly train** the model for a supervised task (e.g., node classification)

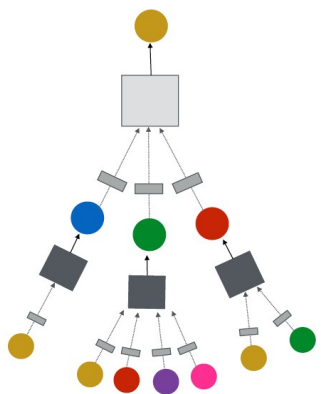


Supervised Training

- **Directly train** the model for a supervised task (e.g., **node classification**)
 - Use cross entropy loss

$$\mathcal{L} = - \left(\sum_{v \in V} y_v \log \left(\sigma(\mathbf{z}_v^T \boldsymbol{\theta}) \right) + (1 - y_v) \log \left(1 - \sigma(\mathbf{z}_v^T \boldsymbol{\theta}) \right) \right)$$

Safe or toxic
drug?



Encoder output:
node embedding

Node class
label

Classification
weights

Unsupervised Training

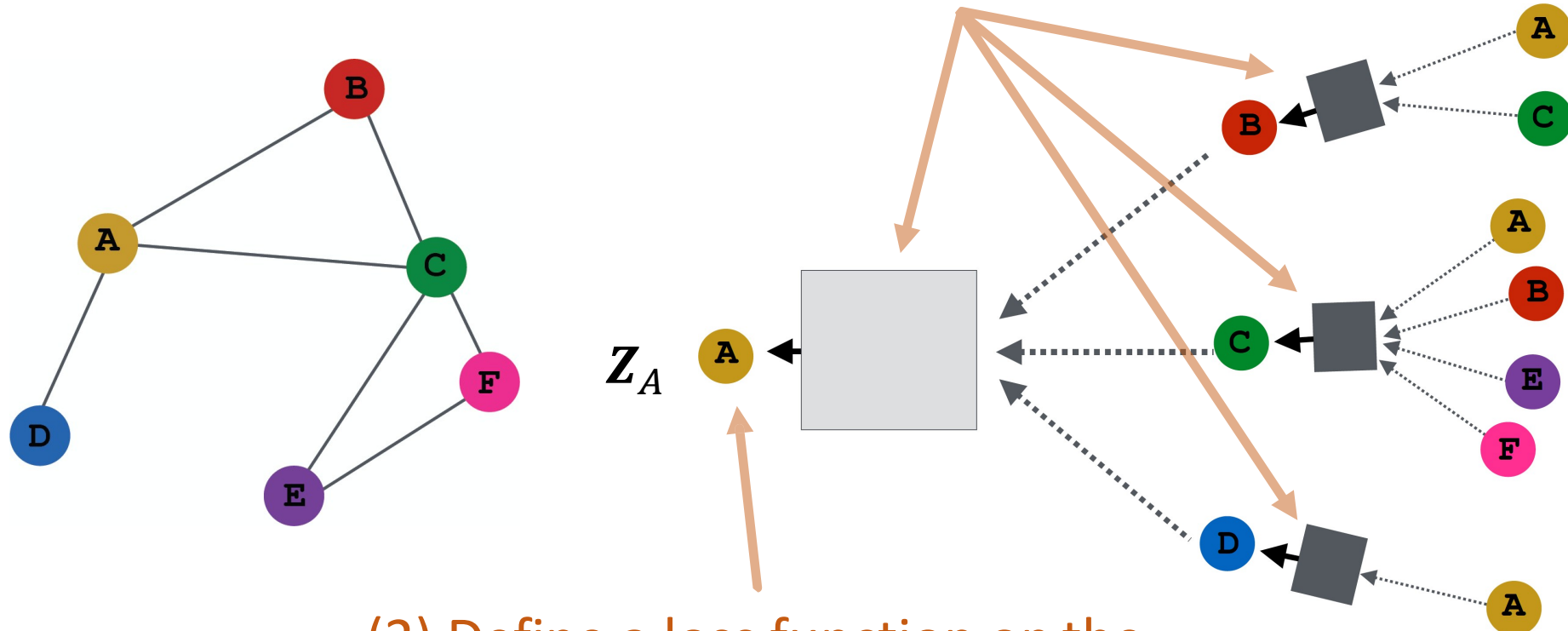
- “Similar” nodes have similar embeddings

$$\mathcal{L} = \sum_{z_u, z_v} \text{CE} \left(y_{u,v}, \text{DEC}(z_u, z_v) \right)$$

- Where $y_{u,v} = 1$ when node u and v are **similar**
- **CE** is the cross entropy
- **DEC** is the decoder such as inner product
- **Node similarity** can be a loss based on:
 - **Random walks** (node2vec, DeepWalk, struc2vec)
 - **Node proximity in the graph**

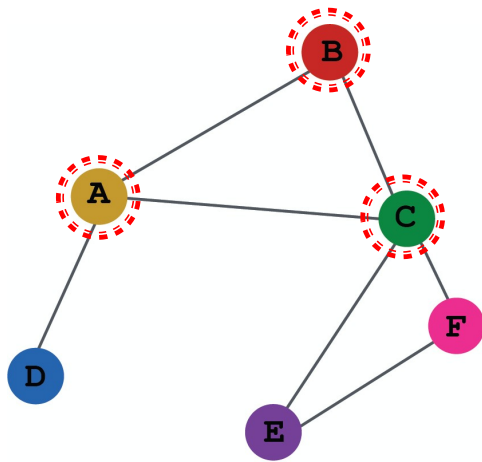
Model Design: Overview

(1) Define a neighborhood aggregation function



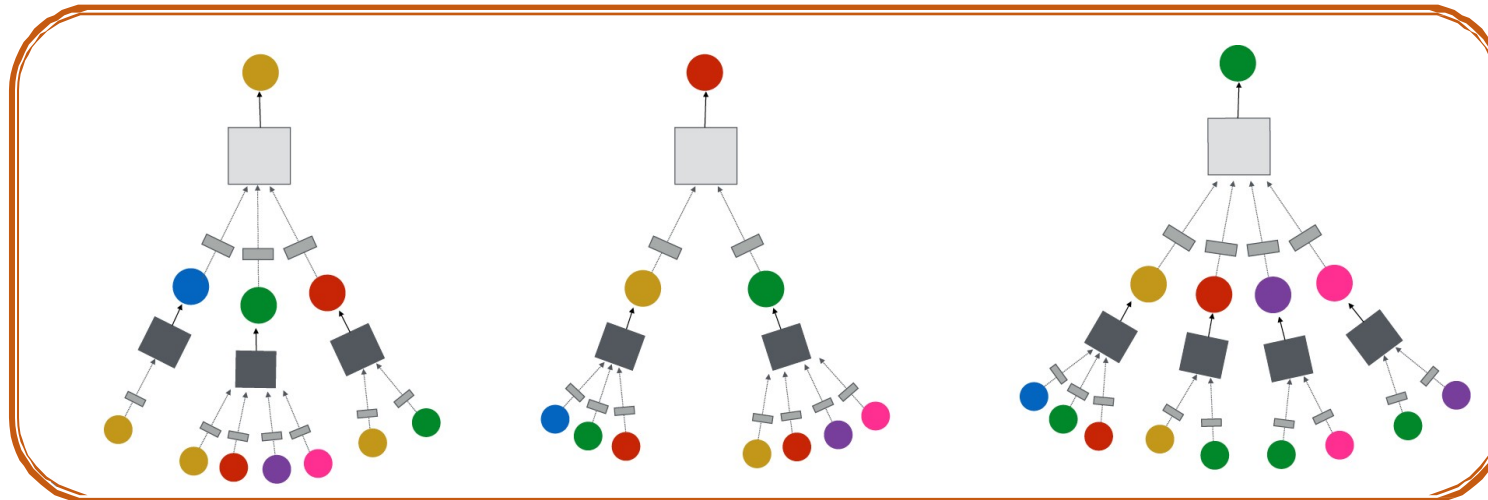
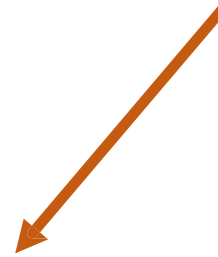
(2) Define a loss function on the embeddings

Model Design: Overview

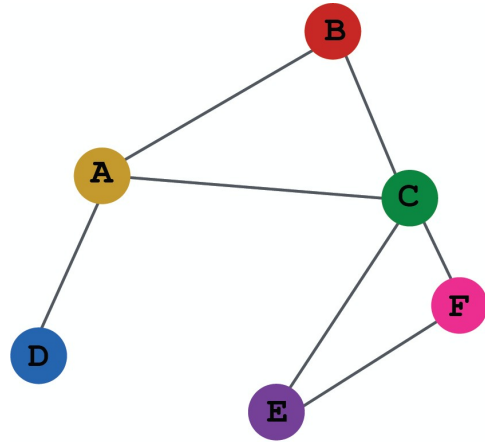


INPUT GRAPH

(3) Train on a set of nodes, i.e.,
a batch of compute graphs



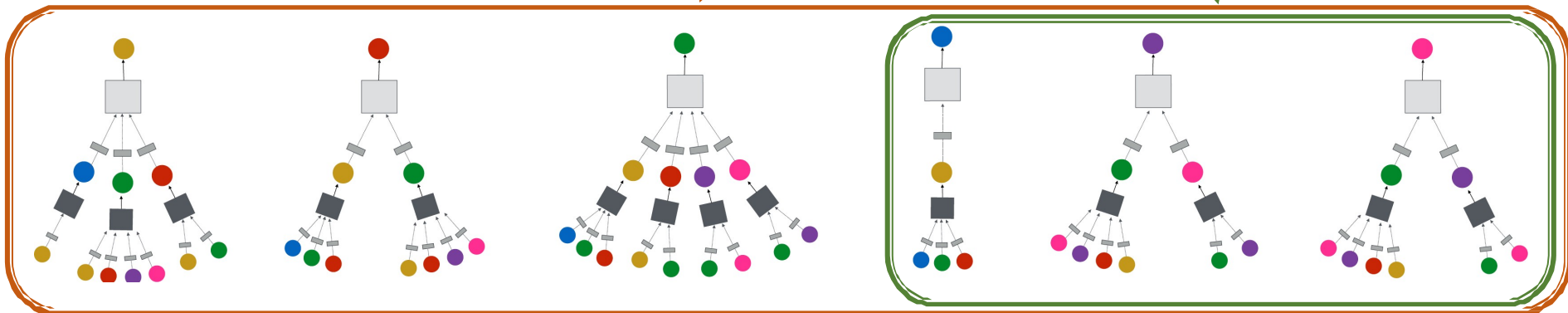
Model Design: Overview



INPUT GRAPH

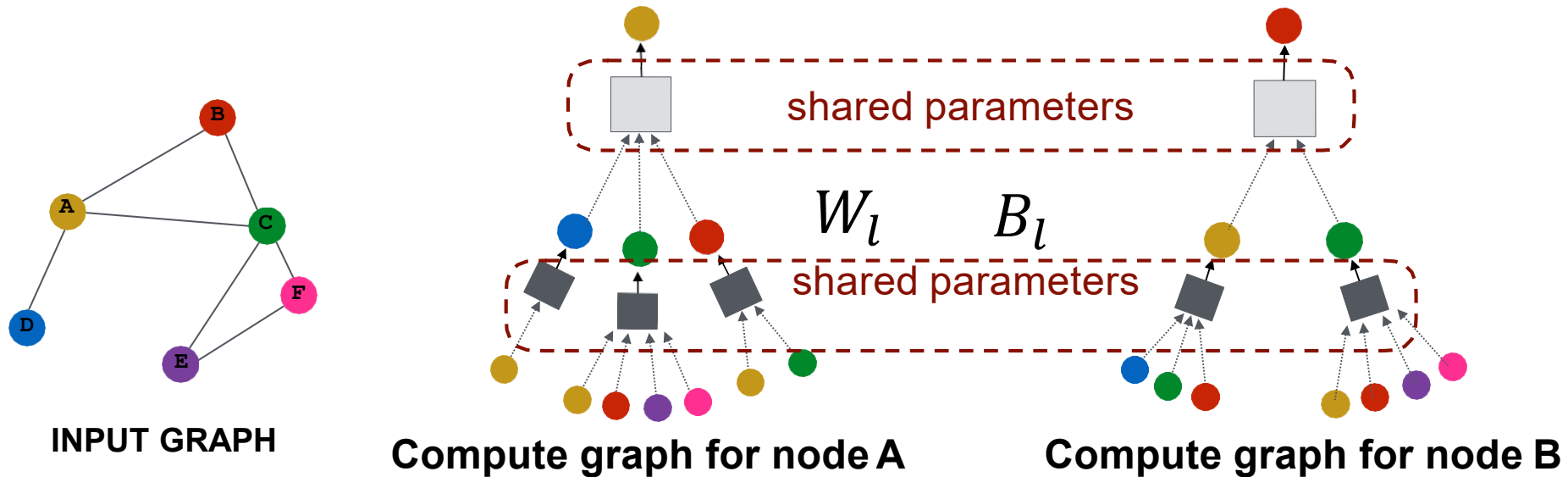
**(4) Generate embeddings
for nodes as needed**

**Even for nodes we never
trained on!**

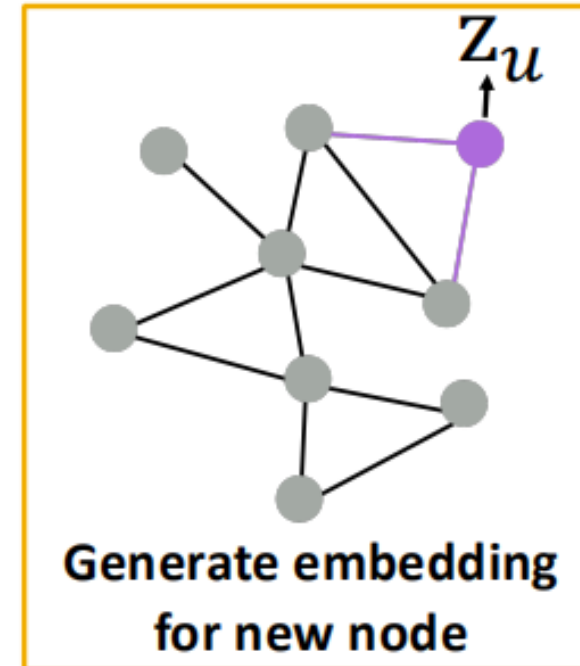
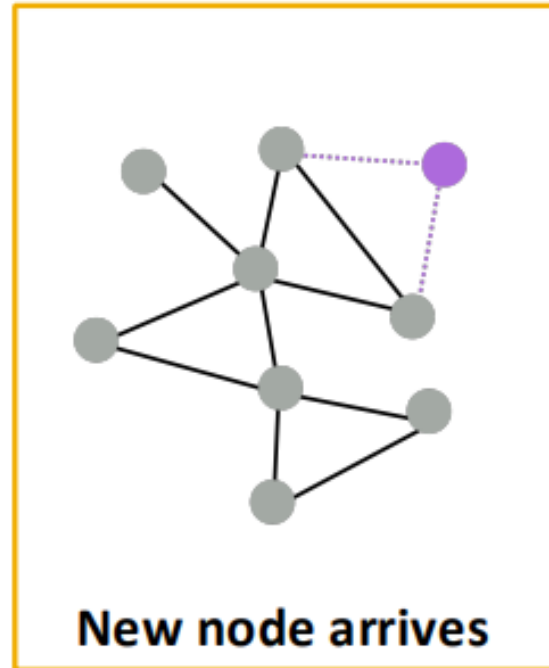
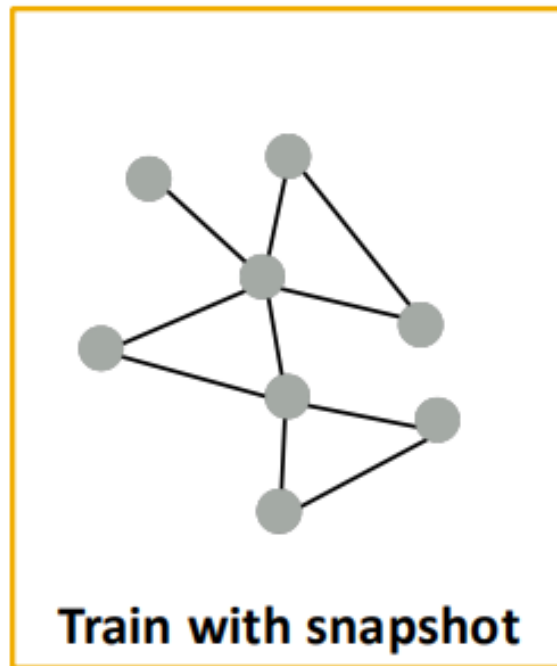


Inductive Capability

- **The same aggregation parameters are shared for all nodes:**
 - The number of model parameters is sublinear in $|V|$ and we can **generalize to unseen nodes!**

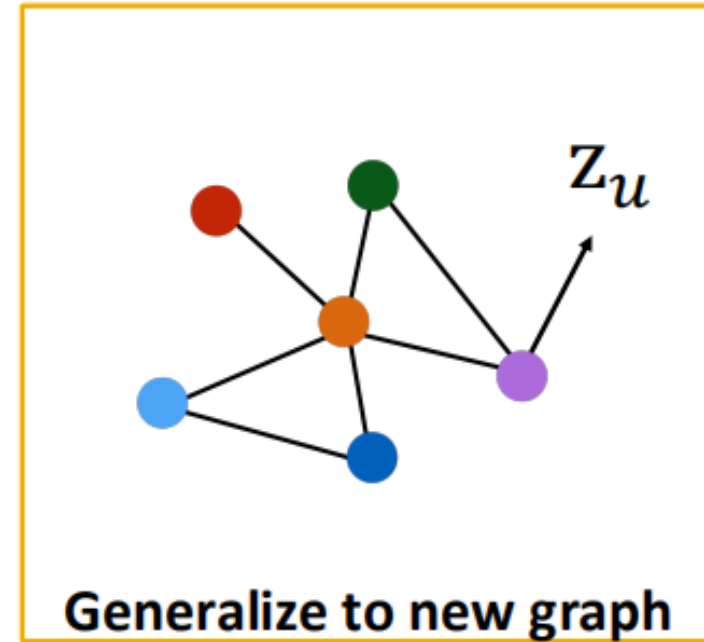
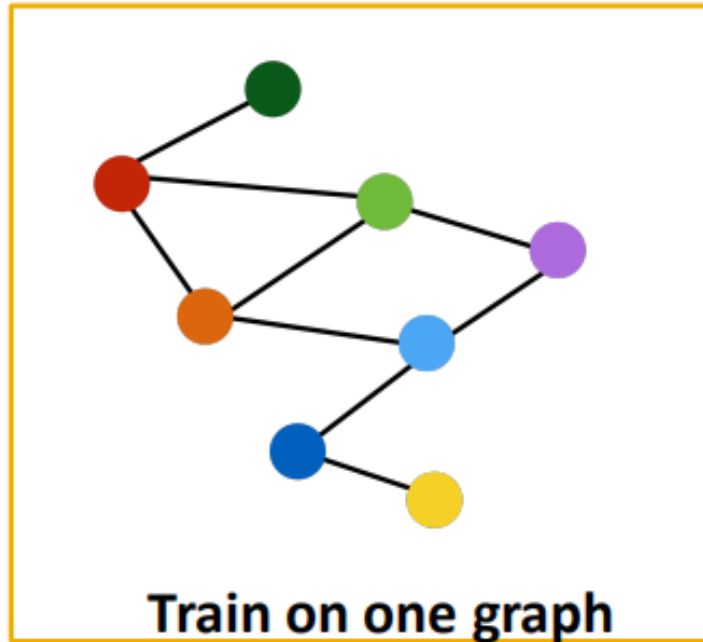


Inductive Capability: New Graphs



- Many application settings constantly encounter previously unseen nodes:
 - E.g., Reddit, YouTube, Google Scholar
- Need to generate new embeddings “on the fly”

Inductive Capability: New Graphs



Inductive node embedding \longrightarrow Generalize to entirely unseen graphs

E.g., train on protein interaction graph from model organism A and generate embeddings on newly collected data about organism B