

Relationship Modeling: Graph Neural Networks (GNNs)

Ju Sun
Computer Science & Engineering

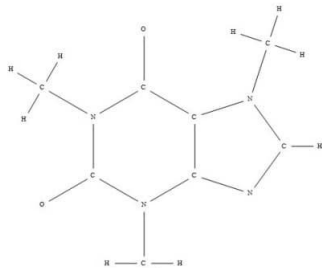
Apr 14, 2026

Outline

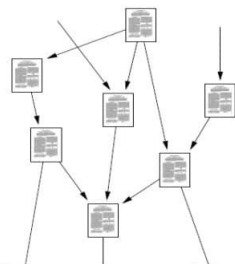
- Why graphs?
- Graphs: basic notions
- Graph neural networks
- Scaling up training

Why graphs?

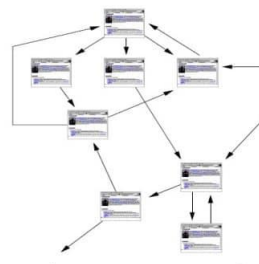
Graphs are everywhere!



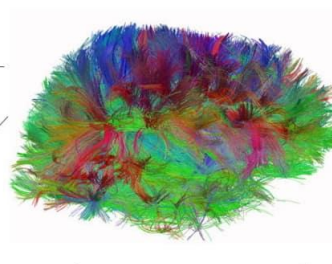
Molecules



Knowledge



Information



Brain/neurons

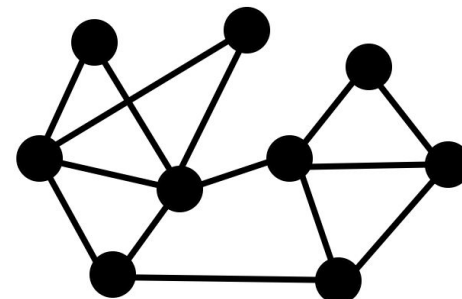
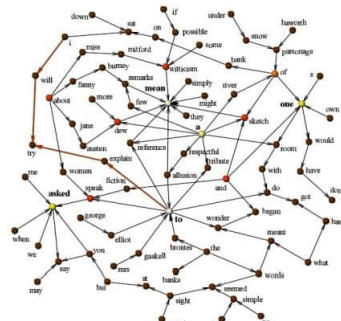


Image credit: Stanford CS224W



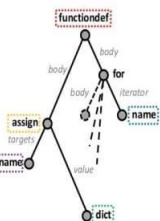
Genes



Communication

```
def encode(obj):  
    """  
    Encode a (possibly nested)  
    dictionary containing complex values  
    into a form that can be serialized  
    using JSON.  
    """  
    if isinstance(obj, dict):  
        return {key: encode(value) for key, value in obj.items()}  
    elif isinstance(obj, list):  
        return [encode(value) for value in obj]  
    elif isinstance(obj, complex):  
        return {'type': 'complex',  
                'r': value.real,  
                'i': value.imag}  
    return obj
```

Software



Social

Graphs model
**relationships/
interactions**

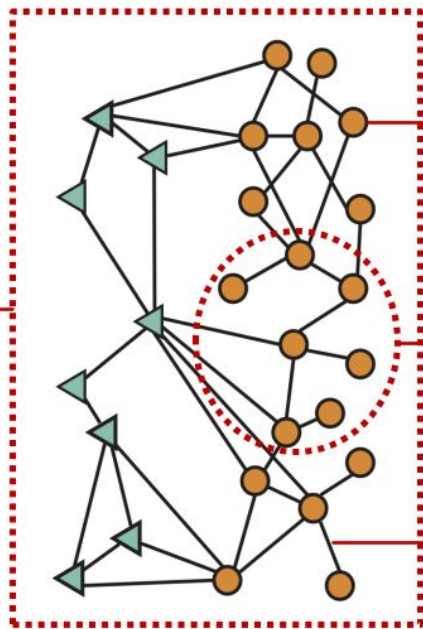
Different tasks on graphs

The whole graph is **part of** all available data, i.e., a data point

Graph prediction (classification/regression), graph generation

A graph is a data point

Graph-level prediction, Graph generation



The whole graph is all available data

Node level

Node prediction (classification/regression)
given labels over part of the graph

Community (subgraph) level

Community detection/clustering
Clustering nodes/edges

Edge-level

Link prediction
Predict links should exists or not

Image credit: Stanford CS224W

Example task 1: Protein folding

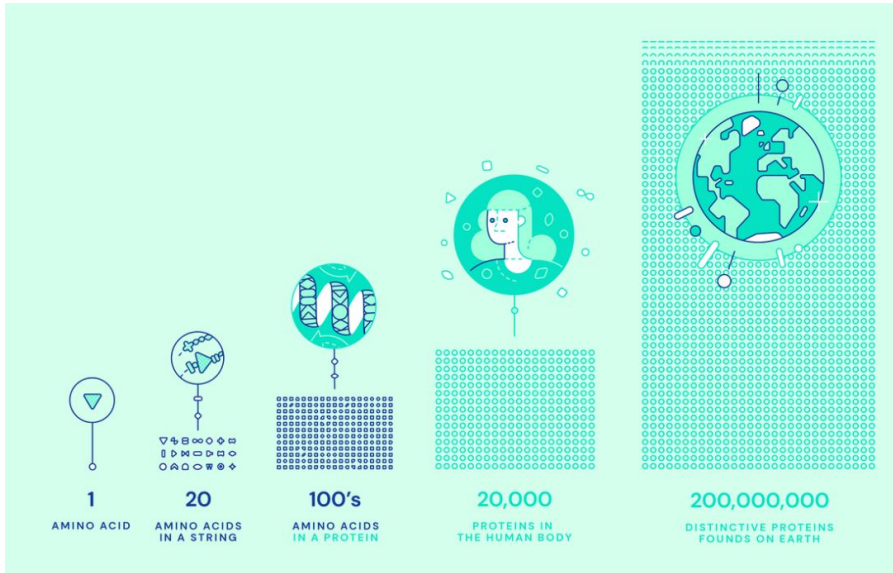
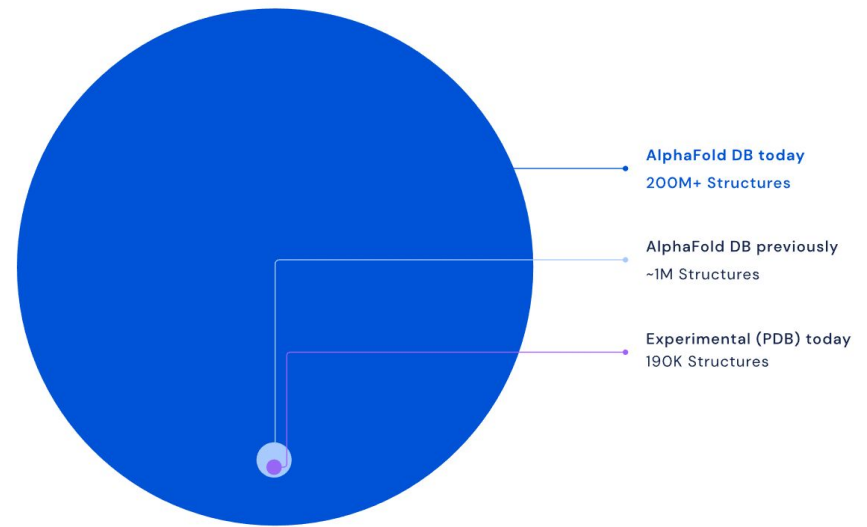
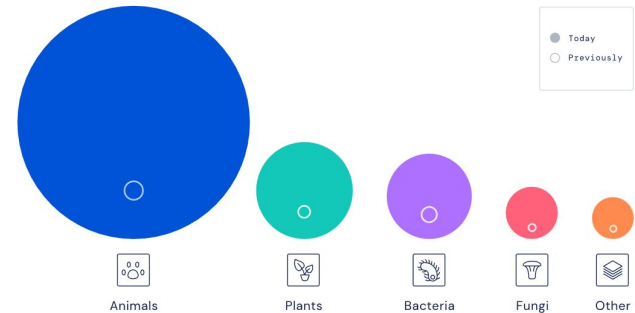


Image credit: Deep Mind



Number of species represented in AlphaFold DB

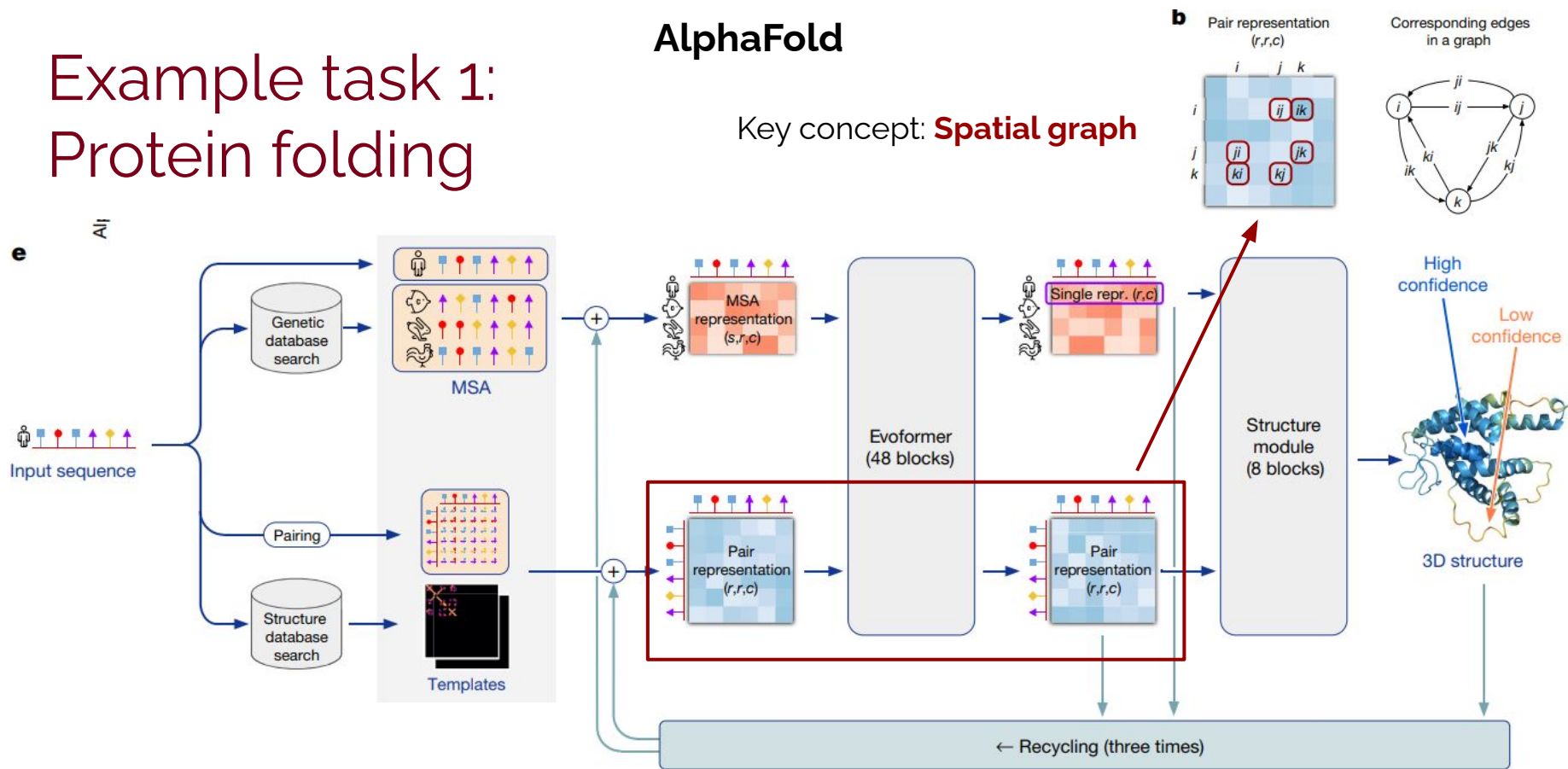
Total increase from ~10K to ~1M



Example task 1: Protein folding

AlphaFold

Key concept: **Spatial graph**



Complete story: <https://www.deepmind.com/research/highlighted-research/alphafold>

Paper: <https://www.nature.com/articles/s41586-021-03819-2>

THE NOBEL PRIZE IN CHEMISTRY 2024



David
Baker

“for computational
protein design”

Demis
Hassabis

“for protein structure prediction”

John M.
Jumper

THE ROYAL SWEDISH ACADEMY OF SCIENCES

nature

Explore content ▾ About the journal ▾ Publish with us ▾ Subscribe

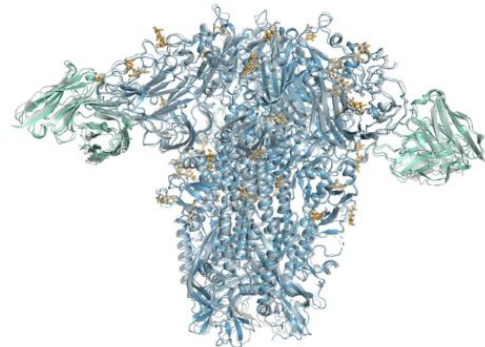
nature > news > article

NEWS | 27 March 2025

AlphaFold is running out of data – so drug firms are building their own version

Thousands of 3D protein structures locked up in big-pharma vaults will be used to create a new AI tool that won't be open to academics.

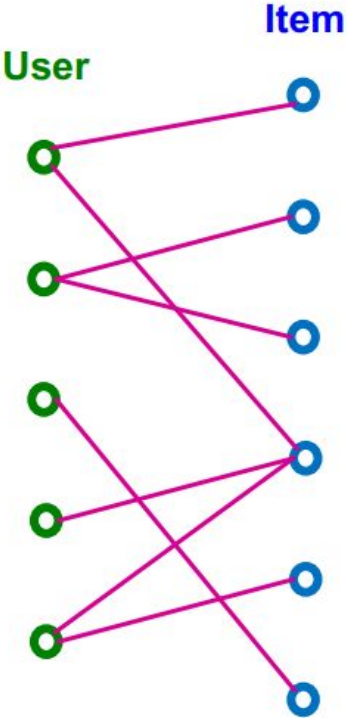
By [Ewen Callaway](#)



An AlphaFold 3 model of a common cold spike protein (blue) interacting with antibodies (green). Credit: Google DeepMind

<https://www.nature.com/articles/d41586-025-00868-9>

Example task 2: Recommendation systems



online shopping,
music/movie
recommendation

Nodes:

Users, items

Edges:

User-item
interactions

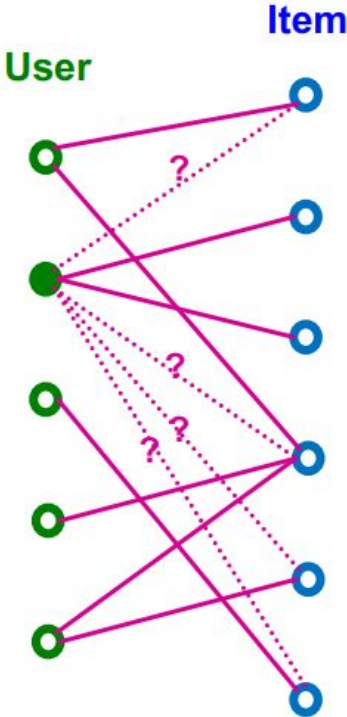
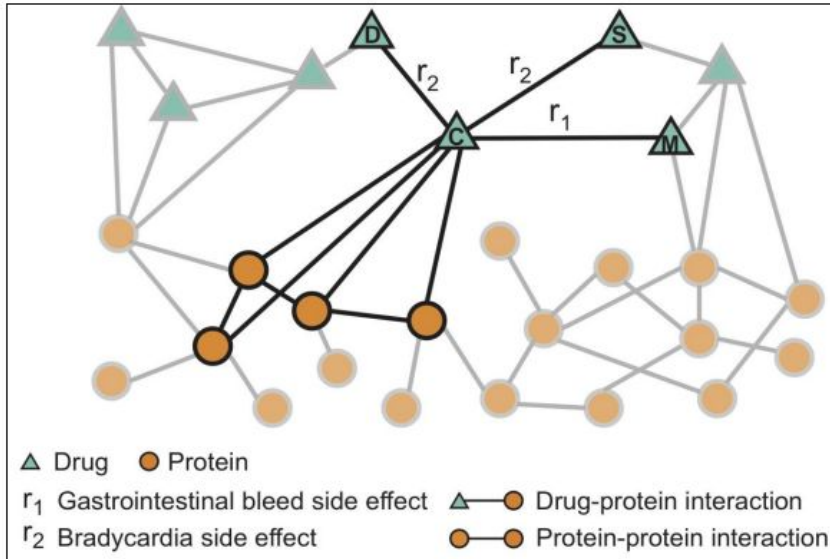


Image credit: Stanford CS224W

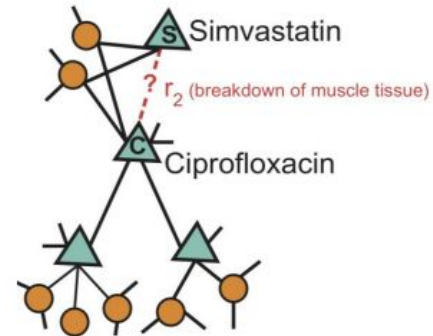
Image credit: Stanford CS224W

Example task 3: Drug adverse effect discovery

- **Nodes:** Drugs & Proteins
- **Edges:** Interactions



Query: How likely will Simvastatin and Ciprofloxacin, when taken together, break down muscle tissue?



Example task 4: Traffic prediction

The screenshot displays a navigation application interface. On the left, there is a search bar with "Your location" and "Minneapolis Institute of Art, 2400 3rd Ave". Below the search bar, there are icons for different travel modes: car, bus, walking, bicycle, and airplane. The selected mode is "Bus".

Below the mode selection, there is a "Leave now" button and an "Options" link. A "Send directions to your phone" button is also present.

The main section shows two transit options:

- Option 1:** 1:50 PM—2:25 PM, 35 min. The route is shown as walking to a bus stop, then taking bus 2, and finally walking. The estimated time is 1:54 PM from Washington Ave & Coffman Union, with a 11 min walk and a 15 min bus ride. A "Details" link is provided.
- Option 2:** 1:49 PM—2:30 PM, 41 min.

On the right, a map of Minneapolis shows the route highlighted in purple. The route starts at the Minneapolis Institute of Art (E & 3rd Ave S) and ends at Washington Ave & Coffman Union. The map includes labels for various neighborhoods like North Loop, Downtown West, and Cedar-Riverside, as well as landmarks like Boom Island Park and the University of Minnesota.

Example task 4: Traffic prediction

- **Nodes:** Road segments
- **Edges:** Connectivity between road segments
- **Prediction:** Time of Arrival (ETA)

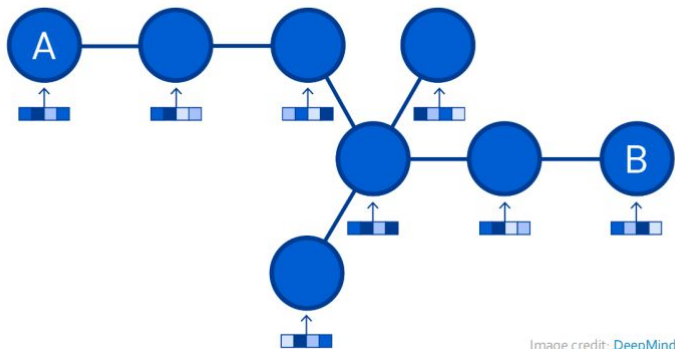


Image credit: Stanford CS224W

Predicting Time of Arrival with Graph Neural Networks

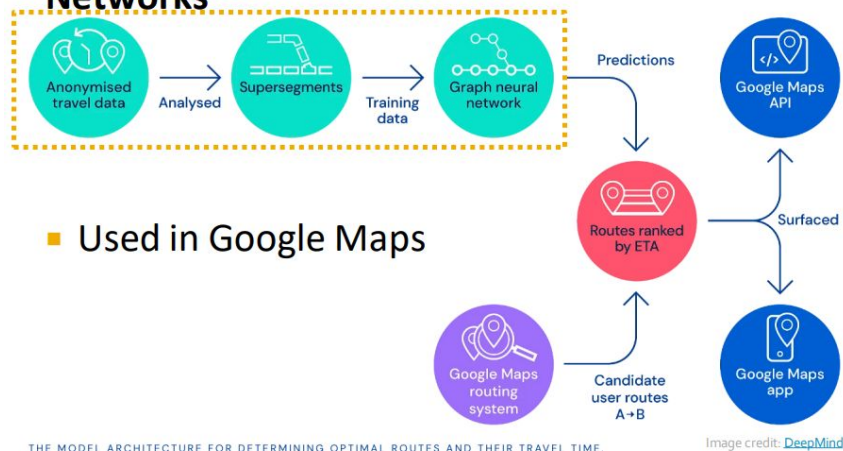


Image credit: Stanford CS224W

Subgraph discovery

Example task 5: Drug discovery

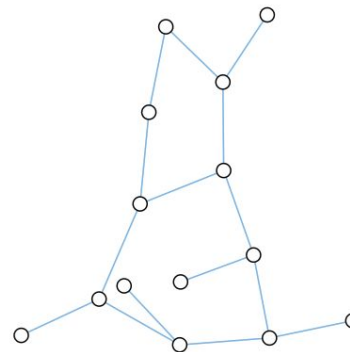
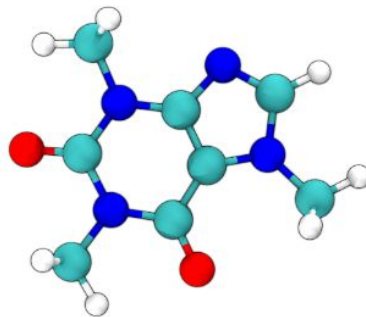
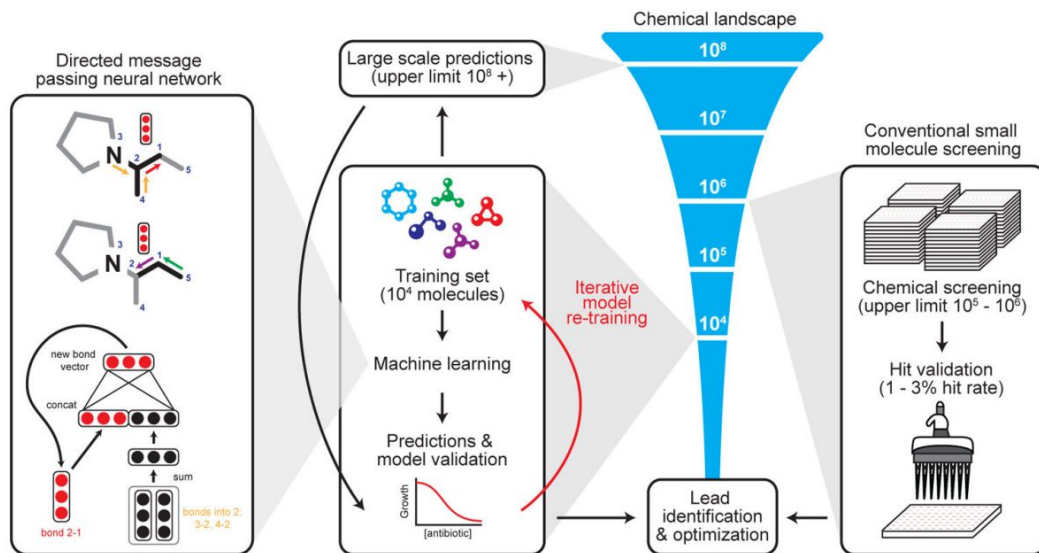


Image credit: <https://distill.pub/2021/gnn-intro/>

Molecules as graphs:

Nodes: atoms

Edges: chemical bounds



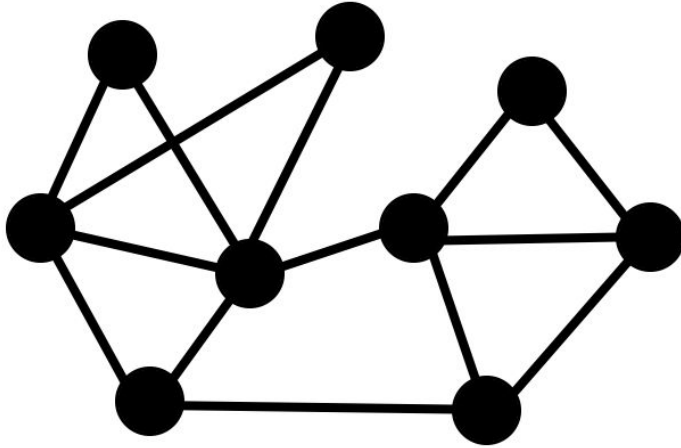
full-graph prediction

Image credit:

<https://pubmed.ncbi.nlm.nih.gov/32084340/>

Graphs: basic notions

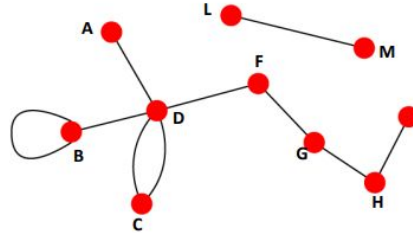
Basic objects



- N : Nodes (also vertices)
- E : Edges (also links)
- $G(N, E)$: Graph

Undirected

- **Links:** undirected (symmetrical, reciprocal)

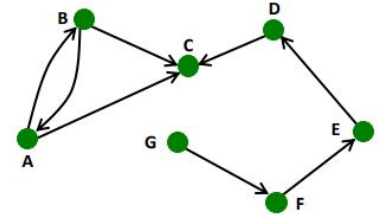


Examples:

- Collaborations
- Friendship on Facebook

Directed

- **Links:** directed (arcs)

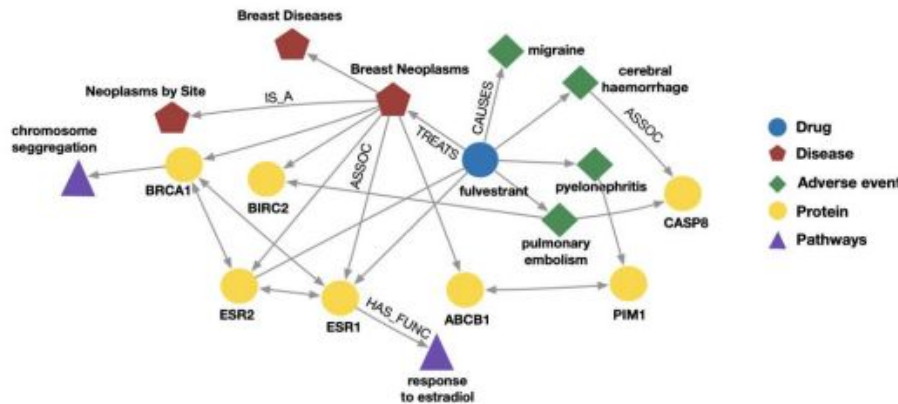


Examples:

- Phone calls
- Following on Twitter

Image credit: Stanford CS224W

Heterogeneous graphs

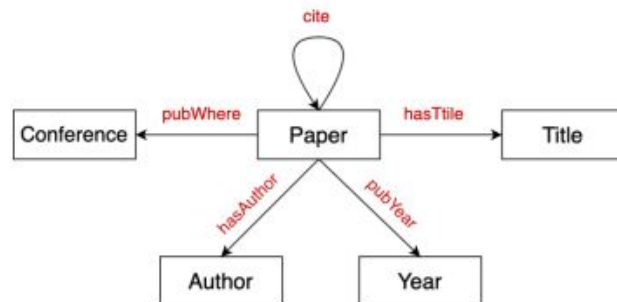


Nodes/Edges are multi-typed

$$G(N, E, T, R)$$

T: types of nodes

R: types of relationships



Biomedical Knowledge Graphs

Example node: Migraine

Example edge: (fulvestrant, Treats, Breast Neoplasms)

Example node type: Protein

Example edge type (relation): Causes

Academic Graphs

Example node: ICML

Example edge: (GraphSAGE, NeurIPS)

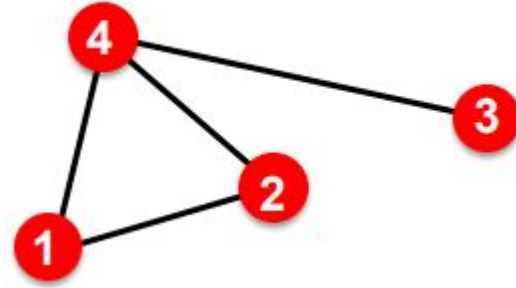
Example node type: Author

Example edge type (relation): pubYear

Image credit: Stanford CS224W

Graph representation

Undirected graphs



$$A = \begin{pmatrix} 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 \end{pmatrix}$$

Adjacency matrix

(1, 4), (1, 2)
(2, 1), (2, 4)
(3, 4)
(4, 1), (4, 2), (4, 3)

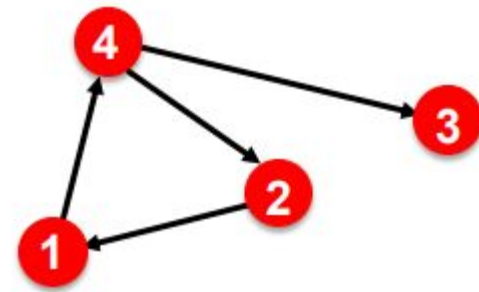
Edge list

1: 2, 4
2: 1, 4
3: 4
4: 1, 2, 3

Adjacency list

Graph representation

Directed graphs



$$A = \begin{pmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 \end{pmatrix}$$

Adjacency matrix

(1, 4)
(2, 1)

(4, 2), (4, 3)

Edge list

1: 4
2: 1
3:
4: 2, 3

Adjacency list

Adjacency matrix is often inefficient



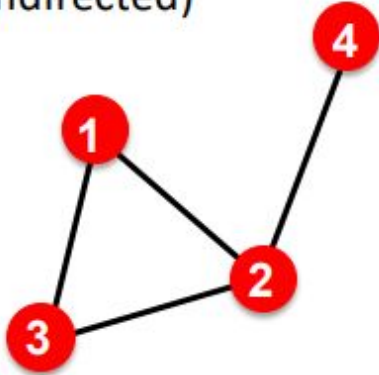
NETWORK	NODES	LINKS	DIRECTED/ UNDIRECTED	N	E
Internet	Routers	Internet connections	Undirected	192,244	609,066
WWW	Webpages	Links	Directed	325,729	1,497,134
Power Grid	Power plants, transformers	Cables	Undirected	4,941	6,594
Phone Calls	Subscribers	Calls	Directed	36,595	91,826
Email	Email Addresses	Emails	Directed	57,194	103,731
Science Collaboration	Scientists	Co-authorship	Undirected	23,133	93,439
Actor Network	Actors	Co-acting	Undirected	702,388	29,397,908
Citation Network	Paper	Citations	Directed	449,673	4,689,479
E. Coli Metabolism	Metabolites	Chemical reactions	Directed	1,039	5,802
Protein Interactions	Proteins	Binding interactions	Undirected	2,018	2,930

$$\text{Density} = |E|/|N|^2$$

Real-world graphs are often very sparse

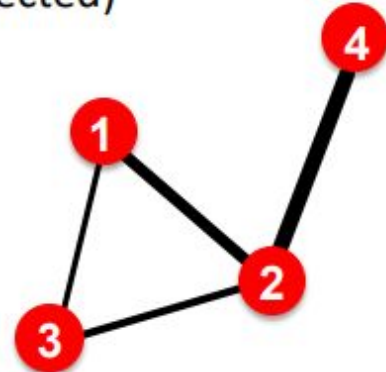
Weighted graphs

- **Unweighted**
(undirected)



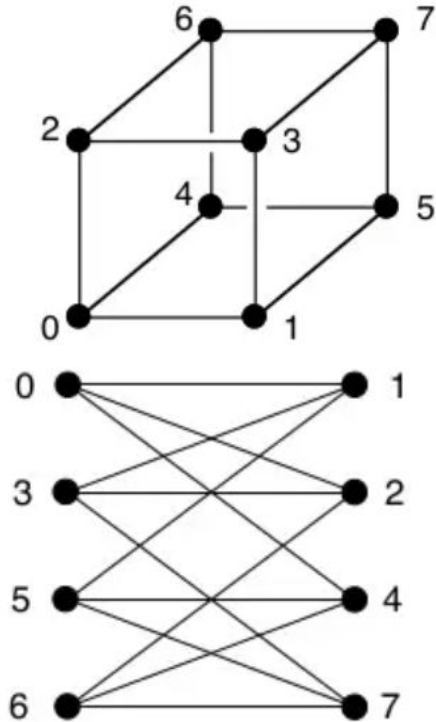
$$\begin{pmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}$$

- **Weighted**
(undirected)



$$\begin{pmatrix} 0 & 2 & 0.5 & 0 \\ 2 & 0 & 1 & 4 \\ 0.5 & 1 & 0 & 0 \\ 0 & 4 & 0 & 0 \end{pmatrix}$$

Graph isomorphism/equivalence



Graph G	Graph H	An isomorphism between G and H
		$f(a) = 1$ $f(b) = 6$ $f(c) = 8$ $f(d) = 3$ $f(g) = 5$ $f(h) = 2$ $f(i) = 4$ $f(j) = 7$

Image credit: https://en.wikipedia.org/wiki/Graph_isomorphism

Isomorphism: there exists a bijective mapping, i.e., **permutation**, results in the same neighborhood structure

Image credit: <https://tonicanada.medium.com/brute-force-code-for-isomorphisms-1241ef180570>

Permutation invariance

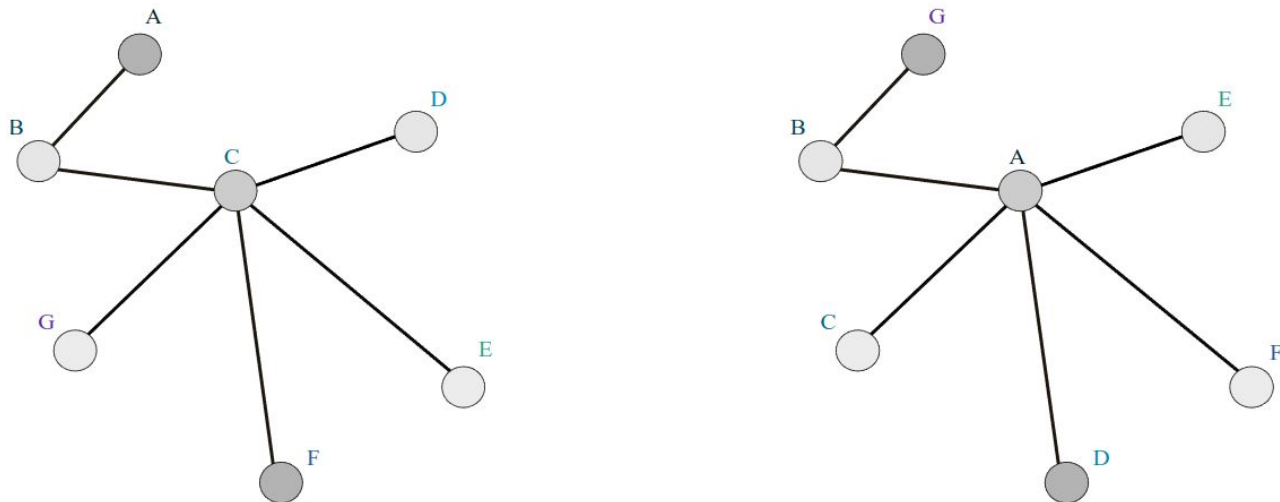


Image credit: Image credit: <https://distill.pub/2021/understanding-gnns/>

Permutation invariance: permuting the names of the nodes doesn't change the graph, as graph nodes are intrinsically orderless

Mathematically: if A is (the adjacency matrix of) a graph, $\Pi A \Pi^T$ is (the adjacency matrix of) an equivalent graph for **any permutation matrix Π**

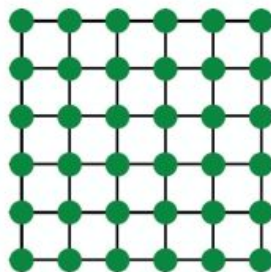
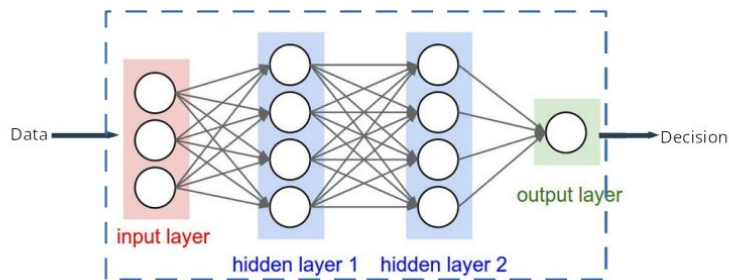
Graph neural networks (GNNs)

Representation learning for graphs

traditional learning pipeline



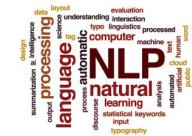
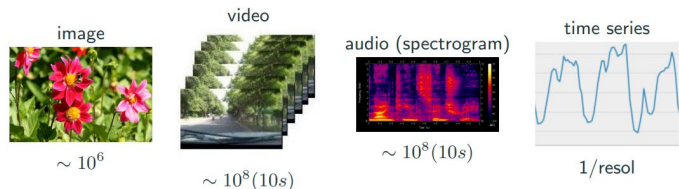
modern learning pipeline



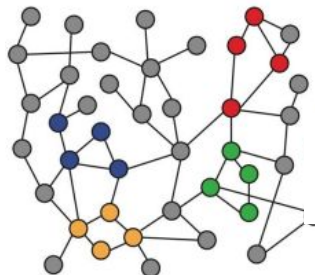
Grid



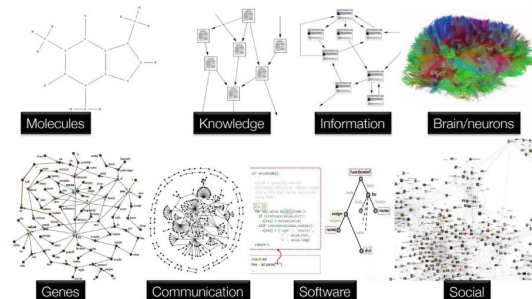
List



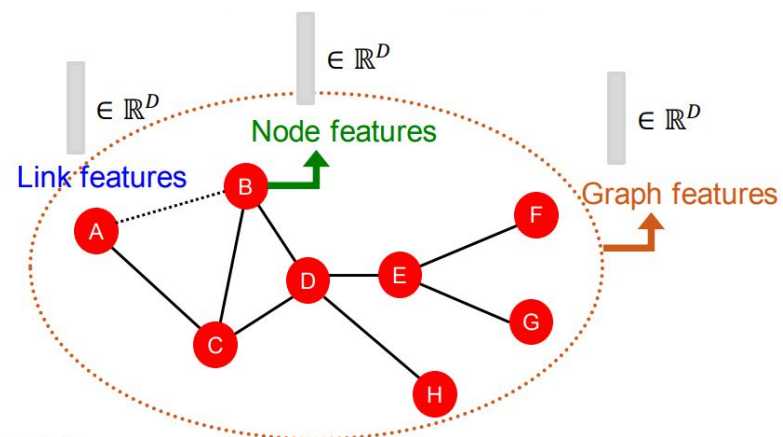
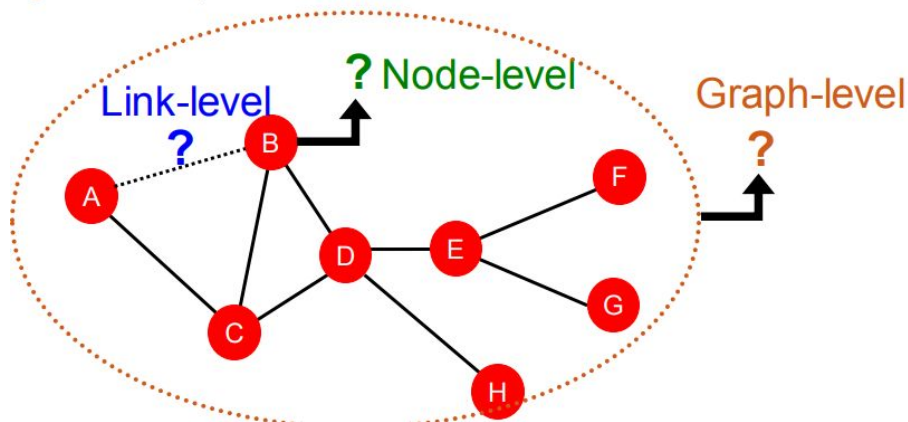
- machine translation, e.g., English \leftrightarrow Chinese
- typing/writing prediction (smart compose)
- semantic classification



Graph



Where to put the features?



- **Train an ML model:**
 - Random forest
 - SVM
 - Neural network, etc.
- **Apply the model:**
 - Given a new node/link/graph, obtain its features and make a prediction

$$\begin{array}{ccc} x_1 & \longrightarrow & y_1 \\ & \vdots & \\ x_N & \longrightarrow & y_N \end{array}$$

$$x \longrightarrow y$$

Node embedding

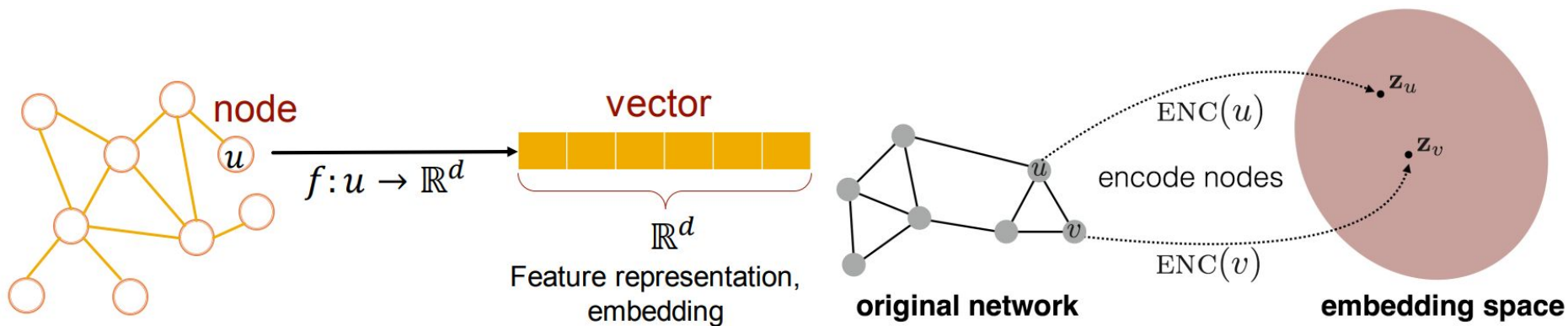


Image credit: Stanford CS224W

N : set of nodes

A : adjacency matrix

$X \in \mathbb{R}^{|N| \times d}$: node (raw) features

u, v : nodes in N

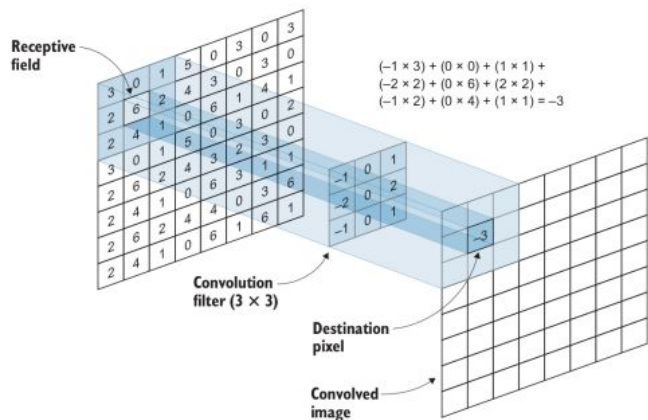
$\mathcal{N}(u)$: neighbors of u

Node raw features: e.g.,

- Biomedical graphs: patient's EHR
- Social network graphs: user profile and images
- When no features: node indicator vector, constant vector

How to define the f ?

We'll bypass fully connected networks directly



(Credit: [Elgandy, 2020])

https://github.com/vdumoulin/conv_arithmetic

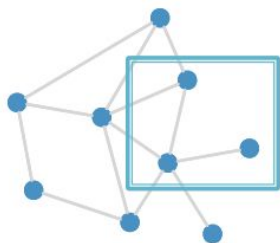
Convolution as performing **local info aggregation** (or **message passing**):

- Each time, the conv window focus on a **local neighborhood** of the current pixel
- Conv effectively **aggregates** the local info by **weighted summation**

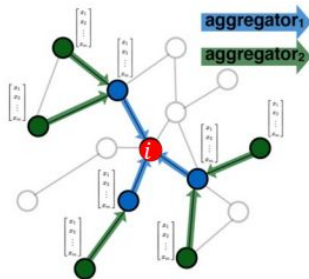
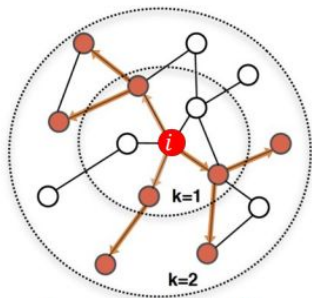
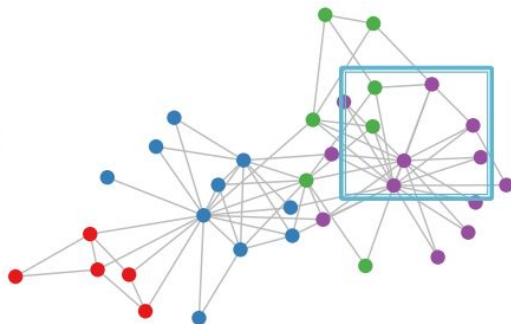
How to define the f ?

Convolution as performing **local info aggregation** (or **message passing**):

- **Neighborhood**: Each time, the conv window focus on a **local neighborhood** of the current pixel
- **Aggregation**: Conv effectively **aggregates** the local info by **weighted summation**



or this:

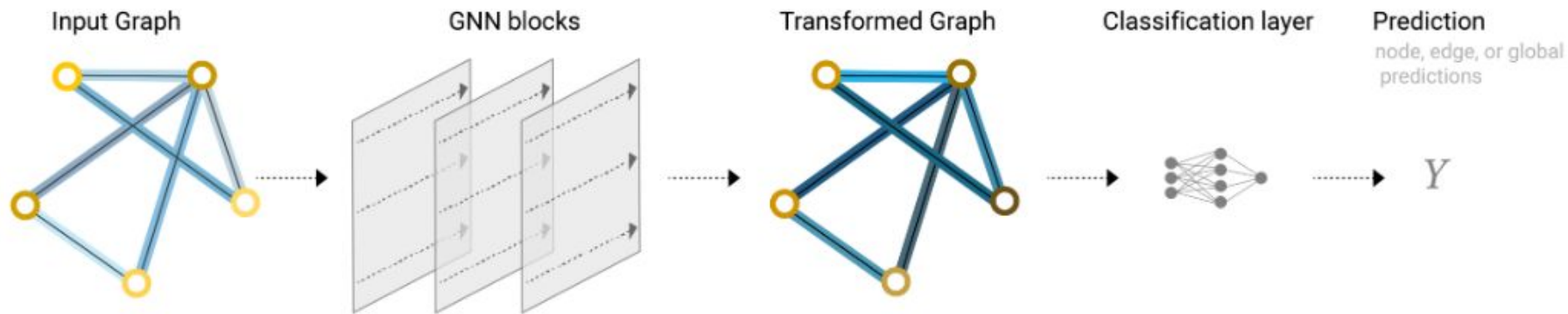


We need neighbors on the graph!

The rigid, grid-based neighborhood doesn't work!

One layer of a graph neural network (GNN)!

Graph in, graph out

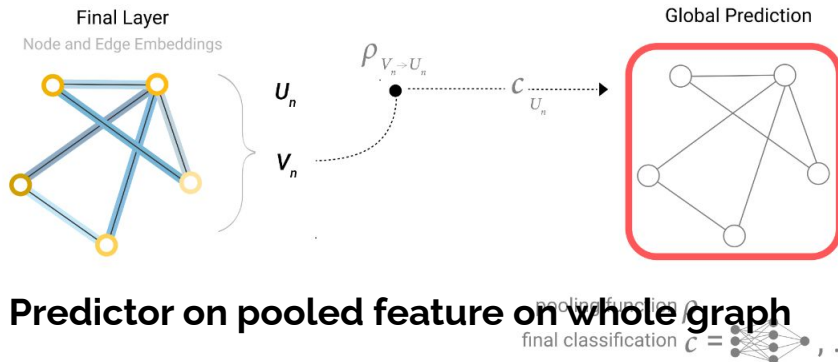
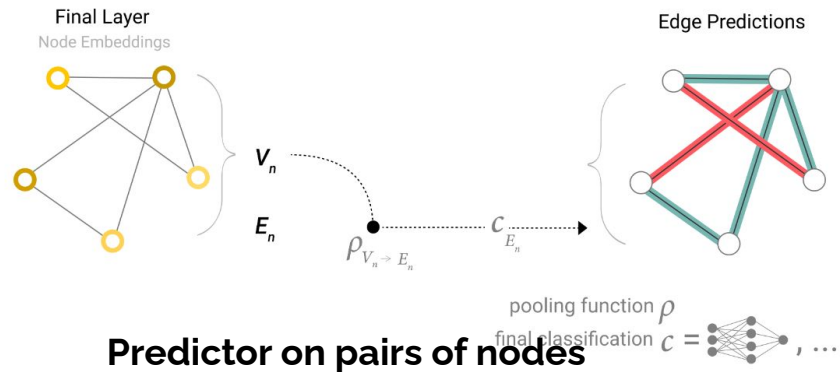
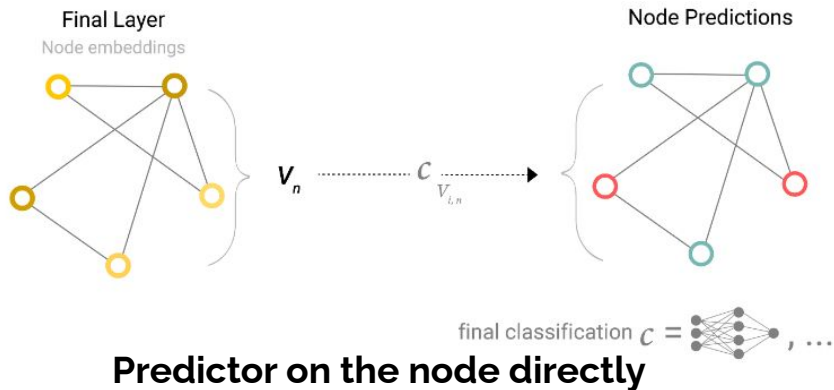


An end-to-end prediction task with a GNN model.

Image credit: <https://distill.pub/2021/gnn-intro/>

How to make supervised predictions?

Image credit: <https://distill.pub/2021/gnn-intro/>

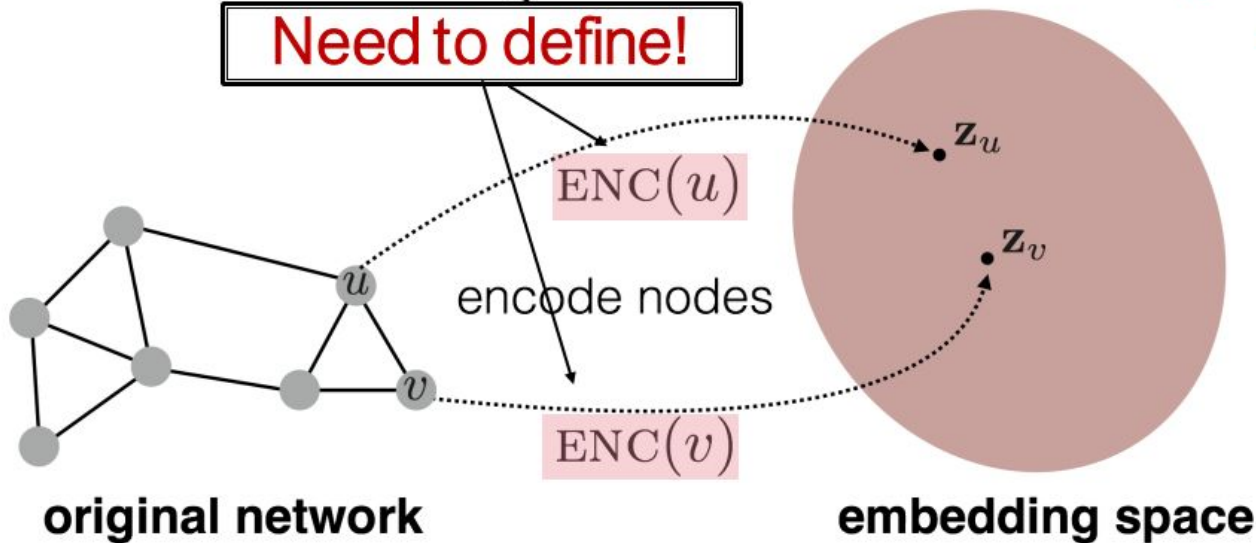


How to perform unsupervised learning?

Goal: $\text{similarity}(u, v)$ $\approx \mathbf{z}_v^T \mathbf{z}_u$
in the original network Similarity of the embedding

Need to define!

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



Look into GNN layers

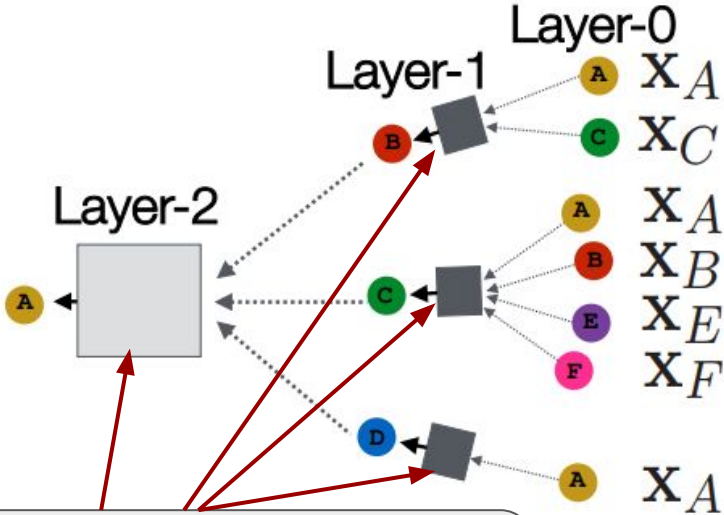
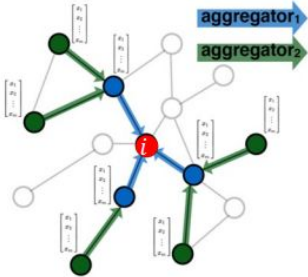
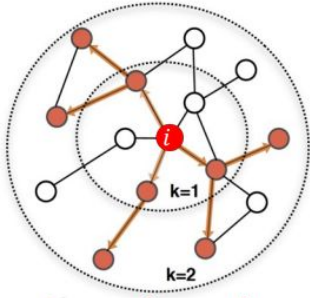
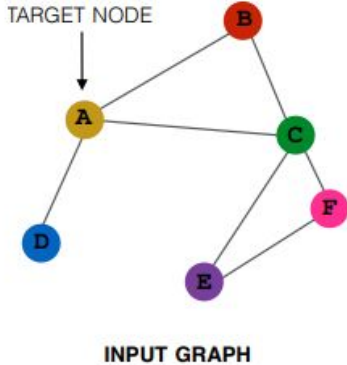


Image credit: Stanford CS224W

These **aggregators** sum up neighboring info, and can be represented by **DNNs**

Permutation invariance

Feature vector for the whole graph

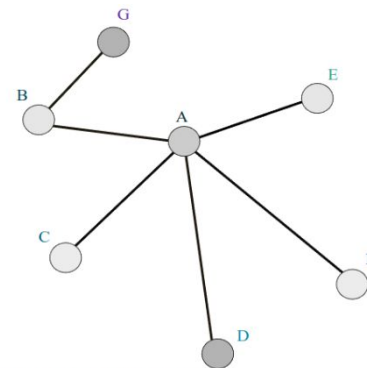
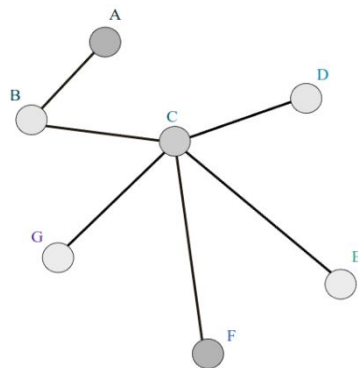
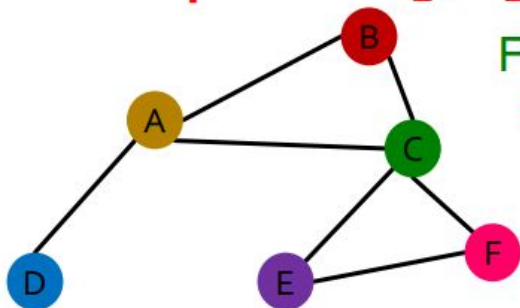


Image credit: <https://distill.pub/2021/understanding-gnns/>

Permutation invariance: permuting the names of the nodes doesn't change the graph, as graph nodes are intrinsically orderless

$$\text{For } f : G(\mathbf{A}, \mathbf{X}) \mapsto \mathbf{h} \cdot f(\mathbf{A}_1, \mathbf{X}_2) = f(\mathbf{A}_2, \mathbf{X}_2)$$

Order plan 1: $\mathbf{A}_1, \mathbf{X}_1$



For two order plans,
output of f should
be the same!

Order plan 2: $\mathbf{A}_2, \mathbf{X}_2$

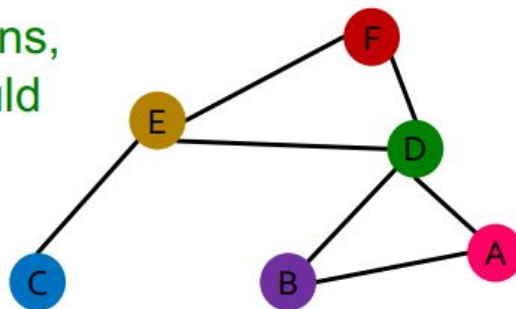


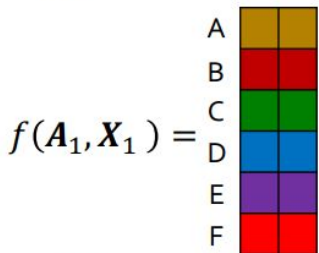
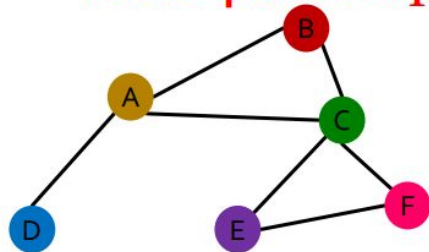
Image credit: Stanford CS224W

Permutation equivariance

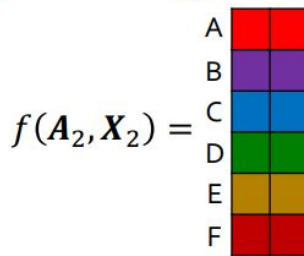
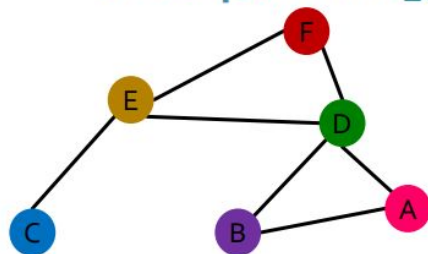
For $f : G(\mathbf{A}, \mathbf{X}) \mapsto \mathbf{H} \in \mathbb{R}^{|N| \times d}$
 $f(\mathbf{\Pi A \Pi}^\top, \mathbf{\Pi X}) = \mathbf{\Pi} f(\mathbf{A}, \mathbf{X})$
 for any permutation $\mathbf{\Pi}$

Collection of feature vectors on all nodes

Order plan 1: $\mathbf{A}_1, \mathbf{X}_1$

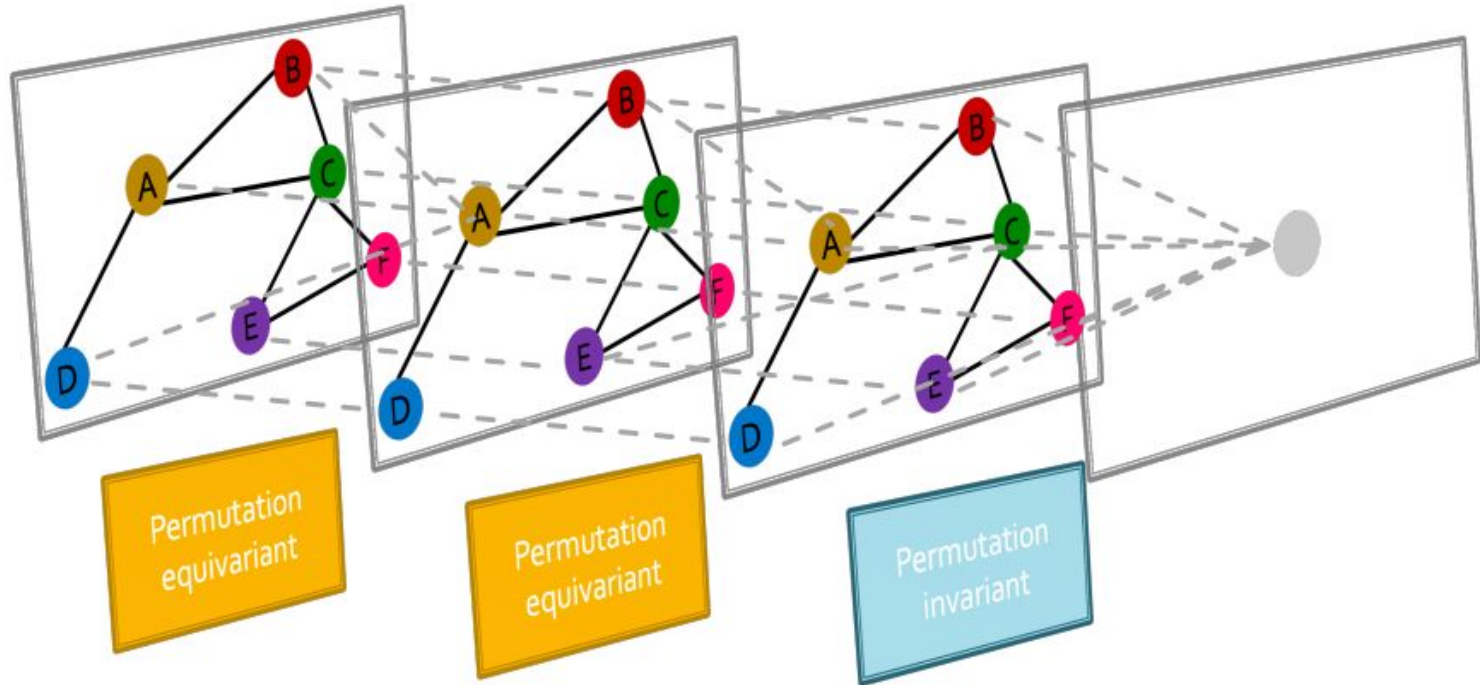


Order plan 2: $\mathbf{A}_2, \mathbf{X}_2$



In other words, if the nodes are re-ordered, the learned features are re-ordered accordingly, **so features are attached to nodes not their names**

A typical GNN consists of multiple permutation equivariant/invariant layers



Graph convolutional networks (GCNs)

$$h_v^{(0)} = x_v \quad \text{for all } v \in V.$$

Node v 's initial embedding. ... is just node v 's original features.

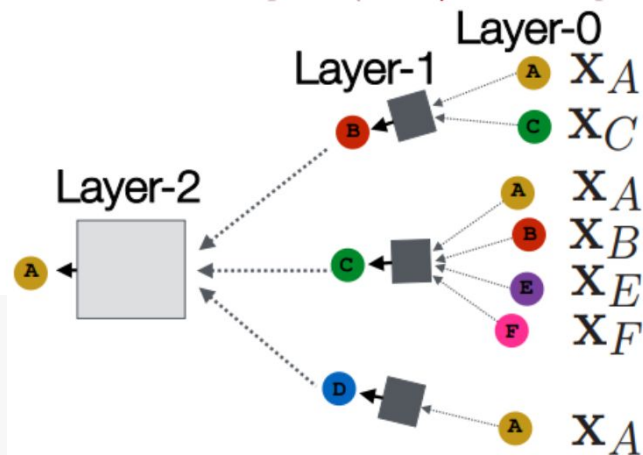
and for $k = 1, 2, \dots$ upto K :

$$h_v^{(k)} = f^{(k)} \left(W^{(k)} \cdot \frac{\sum_{u \in \mathcal{N}(v)} h_u^{(k-1)}}{|\mathcal{N}(v)|} + B^{(k)} \cdot h_v^{(k-1)} \right) \quad \text{for all } v \in V.$$

Node v 's embedding at step k . Mean of v 's neighbour's embeddings at step $k - 1$. Node v 's embedding at step $k - 1$.

Color Codes:

- Embedding of node v .
- Embedding of a neighbour of node v .
- (Potentially) Learnable parameters.



Graph attention networks (GATs)

$$h_v^{(0)} = x_v \quad \text{for all } v \in V.$$

Node v 's initial embedding. ... is just node v 's original features.

and for $k = 1, 2, \dots$ upto K :

$$h_v^{(k)} = f^{(k)} \left(W^{(k)} \cdot \left[\sum_{u \in \mathcal{N}(v)} \alpha_{vu}^{(k-1)} h_u^{(k-1)} + \alpha_{vv}^{(k-1)} h_v^{(k-1)} \right] \right) \quad \text{for all } v \in V.$$

Node v 's embedding at step k .

Weighted mean of v 's neighbour's embeddings at step $k - 1$.

Node v 's embedding at step $k - 1$.

where the attention weights $\alpha^{(k)}$ are generated by an attention mechanism $A^{(k)}$, normalized such that the sum over all neighbours of each node v is 1:

$$\alpha_{vu}^{(k)} = \frac{A^{(k)}(h_v^{(k)}, h_u^{(k)})}{\sum_{w \in \mathcal{N}(v)} A^{(k)}(h_v^{(k)}, h_w^{(k)})} \quad \text{for all } (v, u) \in E.$$

Color Codes:

- Embedding of node v .
- Embedding of a neighbour of node v .
- (Potentially) Learnable parameters.

Graph sample and aggregate (GraphSAGE)

$$h_v^{(0)} = x_v \quad \text{for all } v \in V.$$

Node v 's
initial
embedding.

... is just node v 's
original features.

and for $k = 1, 2, \dots$ upto K :

$$h_v^{(k)} = f^{(k)} \left(W^{(k)} \cdot \left[\text{AGG}_{u \in \mathcal{N}(v)}(\{h_u^{(k-1)}\}), h_v^{(k-1)} \right] \right) \quad \text{for all } v \in V.$$

Node v 's
embedding at
step k .

Aggregation of v 's
neighbour's
embeddings at
step $k - 1$...

... Node v 's
embedding at
step $k - 1$.

... concatenated
with ...

Color Codes:

- Embedding of node v .
- Embedding of a neighbour of node v .
- (Potentially) Learnable parameters.

Graph isomorphism networks (GINs)

$$h_v^{(0)} = x_v \quad \text{for all } v \in V.$$

Node v 's
initial
embedding.

... is just node v 's
original features.

and for $k = 1, 2, \dots$ upto K :

$$h_v^{(k)} = f^{(k)} \left(\sum_{u \in \mathcal{N}(v)} h_u^{(k-1)} + (1 + \epsilon^{(k)}) \cdot h_v^{(k-1)} \right) \quad \text{for all } v \in V.$$

Node v 's
embedding at
step k .

Sum of v 's
neighbour's
embeddings at
step $k - 1$.

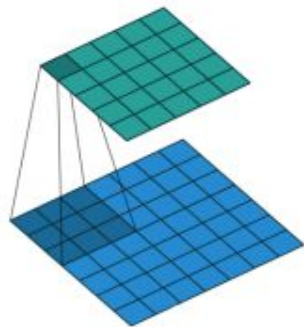
Node v 's
embedding at
step $k - 1$.

Color Codes:

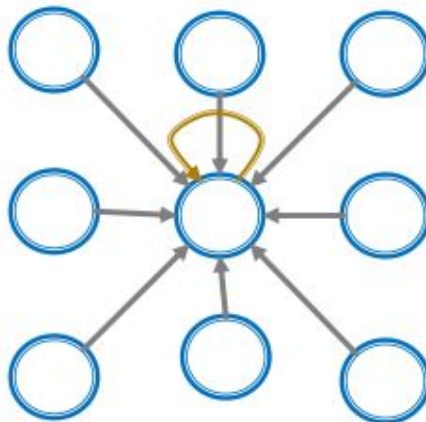
- Embedding of node v .
- Embedding of a neighbour of node v .
- (Potentially) Learnable parameters.

Connection to CNNs and Transformers

filter:



Image



Graph

- CNN is GNN that keeps local ordering
- CNN not permutation-invariant

$$\text{GNN formulation: } h_v^{(l+1)} = \sigma(\mathbf{W}_l \sum_{u \in \mathcal{N}(v)} \frac{h_u^{(l)}}{|\mathcal{N}(v)|} + \mathbf{B}_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$$

$$\text{CNN formulation: } h_v^{(l+1)} = \sigma(\sum_{u \in \mathcal{N}(v)} \mathbf{W}_l^u h_u^{(l)} + \mathbf{B}_l h_v^{(l)}), \forall l \in \{0, \dots, L-1\}$$

Connection to CNNs and Transformers

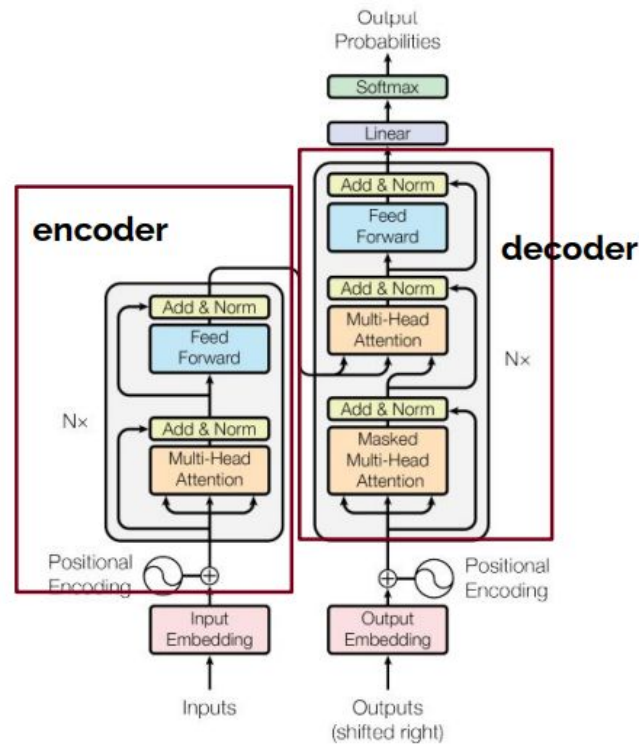
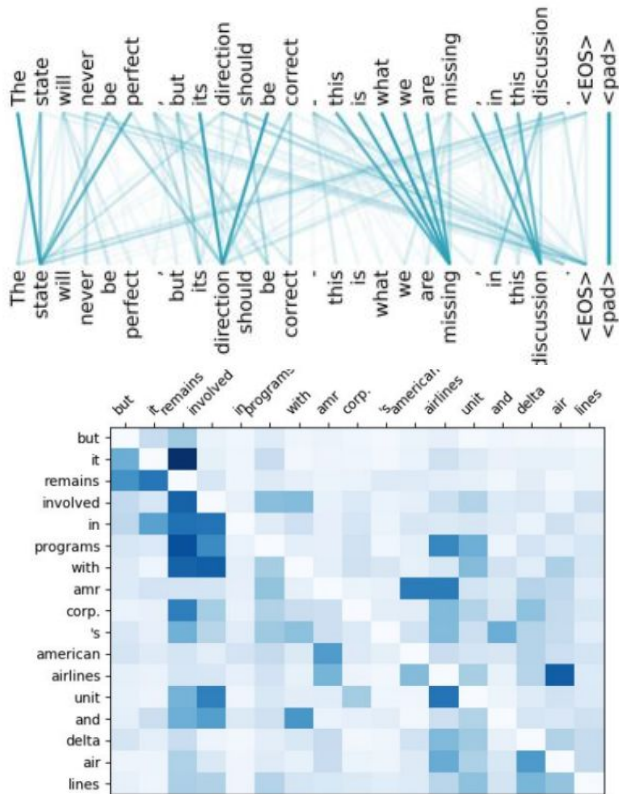


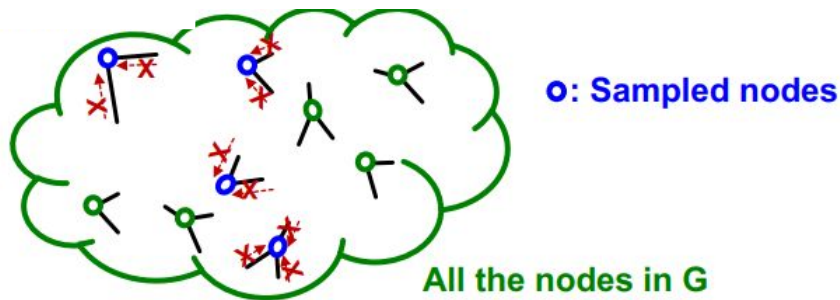
Figure 1: The Transformer - model architecture.

Self-attention (plus feed forward) is a layer of GAT on a complete graph!

Scaling up training

Practical graphs are large yet sparse

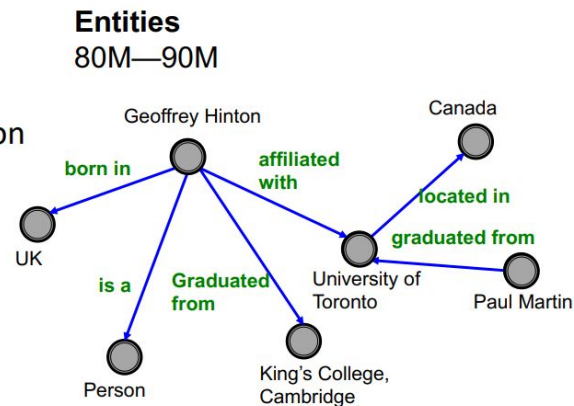
How to perform mini-batch training?



- Mini-batch subsampling induces isolated nodes
- No info to aggregate inside the mini-batch for most nodes

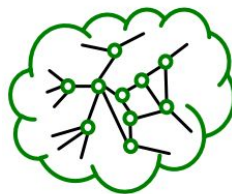
Knowledge Graphs (KGs):

- Wikidata
- Freebase
- **ML tasks:**
 - KG completion
 - Reasoning



Solution: structured subsampling

Large graph



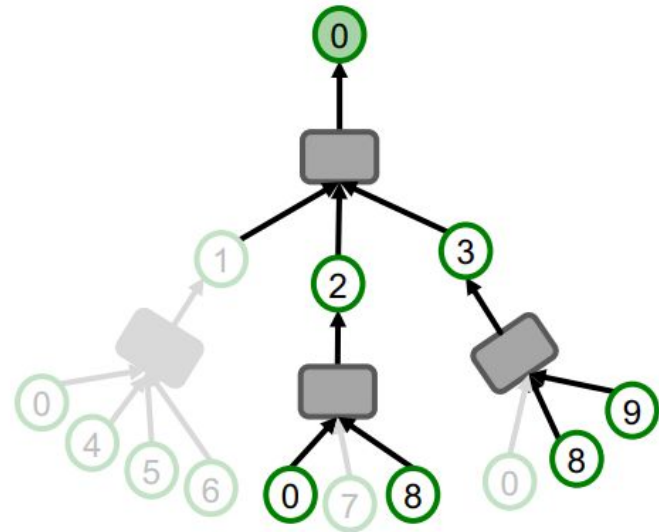
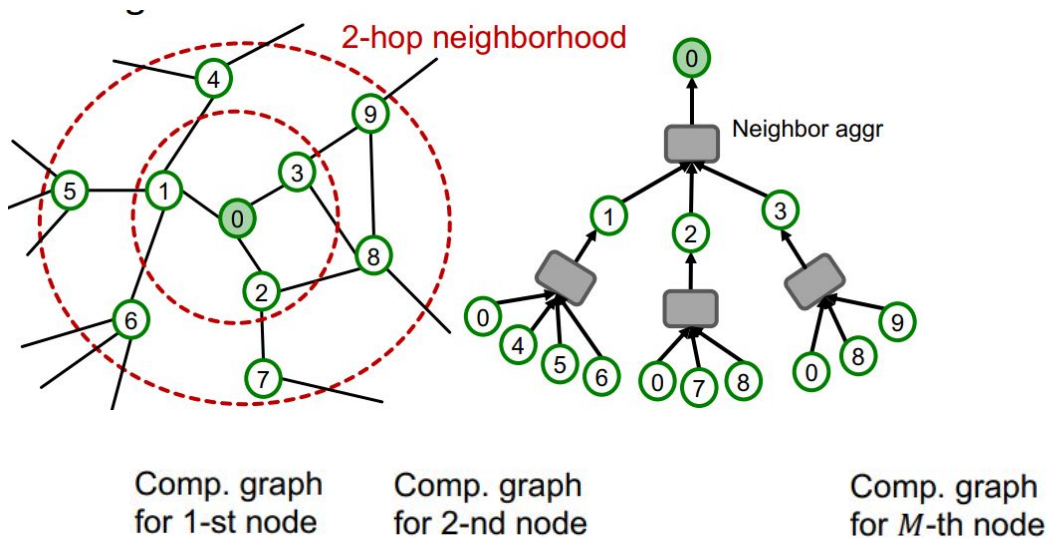
Sampled subgraph
(small enough to be put on a GPU)



Layer-wise
node embeddings
update on the GPU

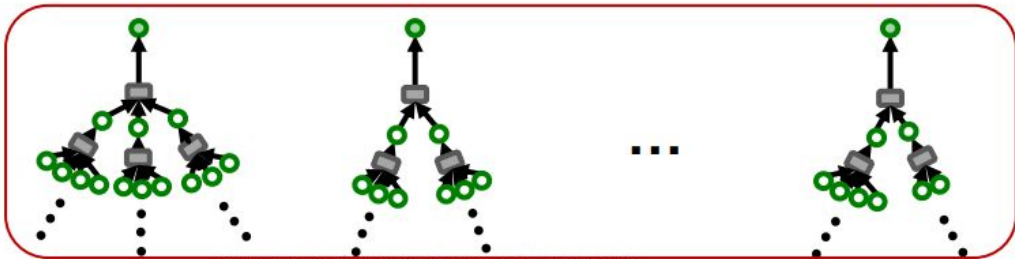


Two structured sub-sampling strategies—I



i.e., **neighborhood sampling**,
instead of iid uniform sampling

Image credit: Stanford CS224W

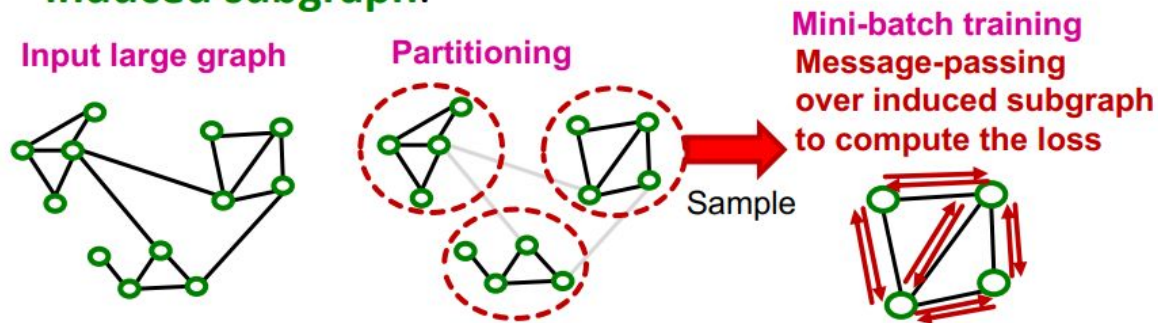


Two structured sub-sampling strategies—II

Cluster-GCN consists of two steps:

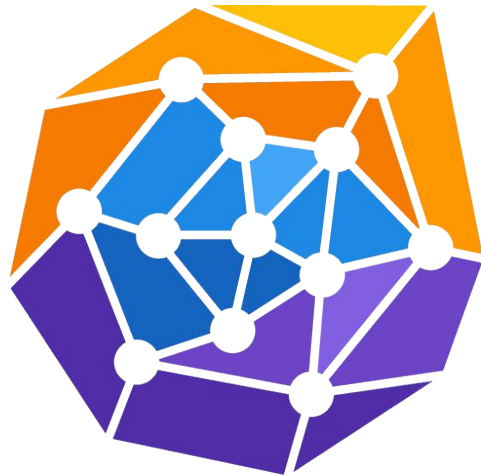
- **Pre-processing**: Given a large graph, partition it into groups of nodes (i.e., subgraphs).
- **Mini-batch training**: Sample one node group at a time. Apply GNN's message passing over the **induced subgraph**.

Rationale: important to keep the community structures, i.e., keep the “backbone” nodes



Software

PyTorch Geometric (PyG)



<https://pytorch-geometric.readthedocs.io/en/latest/>

Deep Graph Library (DGL)



<https://www.dgl.ai/>

Further reading

- What are graph neural networks?
<https://blogs.nvidia.com/blog/2022/10/24/what-are-graph-neural-networks/>
- A Gentle Introduction to Graph Neural Networks
<https://distill.pub/2021/gnn-intro/>
- Understanding Convolutions on Graphs
<https://distill.pub/2021/understanding-gnns/>
- Graph Neural Networks: A Review of Methods and Applications
<https://arxiv.org/abs/1812.08434>
- Stanford CS224W: Machine Learning with Graphs
<https://web.stanford.edu/class/cs224w/index.html>
- Graph Representation Learning https://www.cs.mcgill.ca/~wlh/grl_book/