

# CSCI 566: Deep Learning and its Applications

---

Yue Zhao

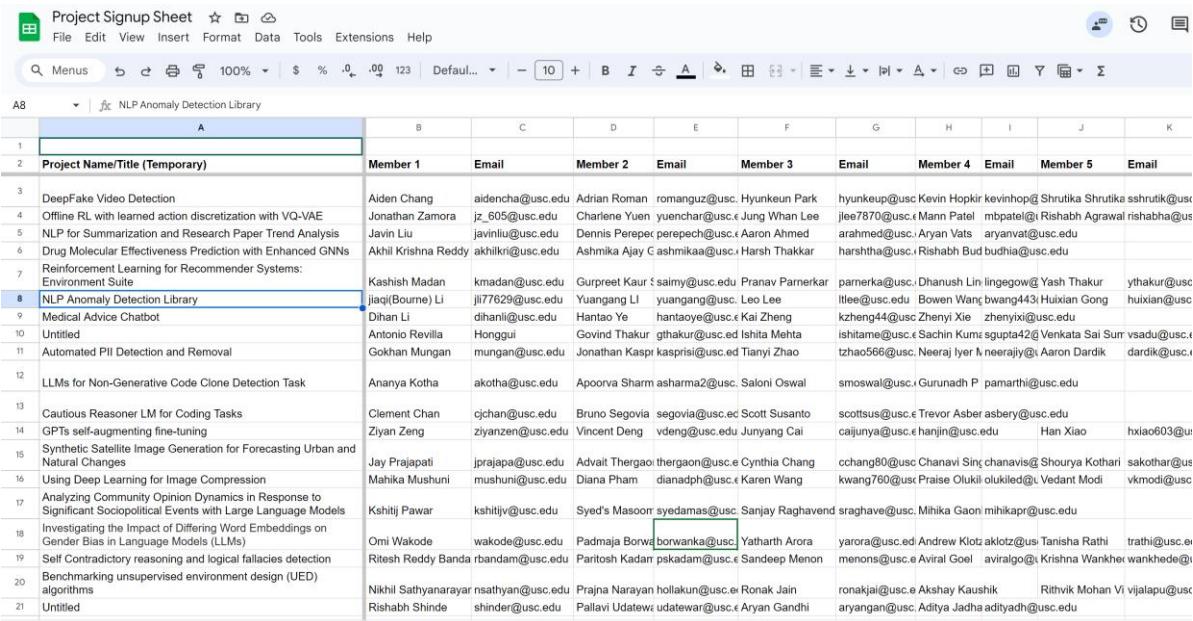
Thomas Lord Department of Computer Science  
University of Southern California

*Credits to previous versions of USC CSCI566,  
CMU 10601/701, Stanford CS 229, 231n, 224w  
Special thanks to Prof. Jure Leskovec on GNN notes.*



# Logistics

Project titles (please add them so we could assign a TA):



The screenshot shows a Google Sheets document with the title "Project Signup Sheet". The table has columns for Project Name/Title (Temporary), Member 1, Email, Member 2, Email, Member 3, Email, Member 4, Email, Member 5, and Email. The rows list various projects and their corresponding team members.

A	B	C	D	E	F	G	H	I	J	K
Project Name/Title (Temporary)	Member 1	Email	Member 2	Email	Member 3	Email	Member 4	Email	Member 5	Email
DeepFake Video Detection	Aiden Chang	aidencha@usc.edu	Adrian Roman	romanguz@usc.edu	Hyunkeun Park	hyunkeup@usc.edu	Kevin Hopkir	kevinhop@usc.edu	Shrutiika Shrutiika	sshrutik@usc.edu
Offline RL with learned action discretization with VQ-VAE	Jonathan Zamora	jz_605@usc.edu	Charlene Yuen	yuerchar@usc.edu	Jung Whan Lee	jlee7870@usc.edu	Mann Patel	mbpatel@usc.edu	Rishabh Agrawal	rishabha@usc.edu
NLP for Summarization and Research Paper Trend Analysis	Javni Liu	javniliu@usc.edu	Dennis Perepec	perepech@usc.edu	Aaron Ahmed	arahmed@usc.edu	Aryan Vats	aryanvat@usc.edu		
Drug Molecular Effectiveness Prediction with Enhanced GNNs	Akhil Krishna Reddy	akhilkri@usc.edu	Ashmika Ajay C	ashmikaa@usc.edu	Harsh Thakkar	harshtha@usc.edu	Rishabh Bud	budhia@usc.edu		
Reinforcement Learning for Recommender Systems: Environment Suite	Kashish Madan	kmadan@usc.edu	Gurpreet Kaur	g.saimy@usc.edu	Pranav Parnerkar	parnerka@usc.edu	Dhanush Lin	lingegow@usc.edu	Yash Thakur	ythakur@usc.edu
NLP Anomaly Detection Library	Jiaqi(Bourne) Li	jli77629@usc.edu	Yuqiang Li	yuqiang@usc.edu	Leo Lee	ltlee@usc.edu	Bowen Wang	bwang443@usc.edu	Huixian Gong	huixian@usc.edu
Medical Advice Chatbot	Dihan Li	dihanli@usc.edu	Hantao Yu	hantaoye@usc.edu	Kai Zheng	kzheng44@usc.edu	Zhenyi Xie	zhenyixi@usc.edu		
Automated PII Detection and Removal	Antonio Revilla	honggui@usc.edu	Govind Thakur	gthakur@usc.edu	Ishita Mehta	ishitame@usc.edu	Sachin Kumar	skupta42@usc.edu	Venkata Sai Sur	vsadu@usc.edu
LLMs for Non-Generative Code Clone Detection Task	Gokhan Mungan	mungan@usc.edu	Jonathan Kaspr	kaspris@usc.edu	Tianyi Zhao	tzhao566@usc.edu	Neeraj Iyer	neerajiy@usc.edu	Aaron Dardik	dardik@usc.edu
Cautious Reasoner LM for Coding Tasks	Ananya Kotha	akotha@usc.edu	Apoorva Sharma	arma2@usc.edu	Saloni Oswal	smoswal@usc.edu	Gurunadh P	pamarthi@usc.edu		
GPTs self-augmenting fine-tuning	Clement Chan	cjchan@usc.edu	Bruno Segovia	segovia@usc.edu	Scott Susanto	scottsus@usc.edu	Trevor Asber	asber@usc.edu		
Synthetic Satellite Image Generation for Forecasting Urban and Natural Changes	Ziyan Zeng	ziyanzen@usc.edu	Vincent Deng	vdeng@usc.edu	Junyang Cai	caijuny@usc.edu	Han Xiao	xhiao603@usc.edu		
Using Deep Learning for Image Compression	Jay Prajapati	jprajapa@usc.edu	Advait Theragoa	theragoa@usc.edu	Cynthia Chang	cchang80@usc.edu	Chanavi Sing	chanavas@usc.edu	Shourya Kothari	sakothar@usc.edu
Analyzing Community Opinion Dynamics in Response to Significant Sociopolitical Events with Large Language Models	Mahika Mushuni	mushuni@usc.edu	Diana Pham	dianaph@usc.edu	Karen Wang	kwang760@usc.edu	Praise Oluklidi	oluklidi@usc.edu	Vedant Modi	vkmodi@usc.edu
Investigating the Impact of Differing Word Embeddings on Gender Bias in Language Models (LLMs)	Kshitij Pawar	kshitijpv@usc.edu	Syed's Masoor	syedmas@usc.edu	Sanjay Raghavend	sraghav@usc.edu	Mihika Gaoni	mihikapr@usc.edu		
Self Contradictory reasoning and logical fallacies detection	Omi Wakode	wakode@usc.edu	Padmaja Borwanka	borwanka@usc.edu	Yatharth Arora	yarora@usc.edu	Andrew Klotz	aklotz@usc.edu	Tanisha Rathi	trathi@usc.edu
Benchmarking unsupervised environment design (UED) algorithms	Ritesh Reddy Banda	rbandam@usc.edu	Paritosh Kadar	pskadam@usc.edu	Sandeep Menon	menons@usc.edu	Aviral Goel	avralgo@usc.edu	Krishna Wankhede	krishna.wankhede@usc.edu
Untitled	Nikhil Sathyaranayana	nsathyan@usc.edu	Prajna Narayan	holakun@usc.edu	Ronak Jain	ronakj@usc.edu	Akshay Kaushik		Rithvik Mohan	viyalapu@usc.edu
	Rishabh Shinde	shinder@usc.edu	Pallavi Udatewar	udatewar@usc.edu	Aryan Gandhi	aryangan@usc.edu	Aditya Jadhav	adityadh@usc.edu		

# Logistics

If you do not put your project here - you consent to do this without TA's help...

Project Signup Sheet

A8 NLP Anomaly Detection Library

	A	B	C	D	E	F	G	H	I	J	K
1	Project Name/Title (Temporary)	Member 1	Email	Member 2	Email	Member 3	Email	Member 4	Email	Member 5	Email
2	DeepFake Video Detection	Aiden Chang	aidencha@usc.edu	Adrian Roman	romanguz@usc.edu	Hyunkeun Park	hyunkeup@usc.edu	Kevin Hopkin	kevinhop@usc.edu	Shrutiika Shrutiika	shrutik@usc.edu
3	Offline RL with learned action discretization with VQ-VAE	Jonathan Zamora	jz_605@usc.edu	Charlene Yuen	yuencchar@usc.edu	Jung Whan Lee	jlee7870@usc.edu	Mann Patel	mpatel@usc.edu	Rishabh Agrawal	rishabha@usc.edu
4	NLP for Summarization and Research Paper Trend Analysis	Javin Liu	javinliu@usc.edu	Dennis Perepec	perepec@usc.edu	Aaron Ahmed	arahmed@usc.edu	Aryan Varaiya	aryanvar@usc.edu		
5	Drug Molecular Effectiveness Prediction with Enhanced GNNs	Akhil Krishna Reddy	akhilkri@usc.edu	Ashmitka Ajay	C ashmitka@usc.edu	Harsh Thakkar	harshtha@usc.edu	Rishabh Budiu	budia@usc.edu		
6	Reinforcement Learning for Recommender Systems: Environment Suite	Kashish Madan	kmadan@usc.edu	Gurpreet Kaur	saimy@usc.edu	Pranav Palmerkar	parnerka@usc.edu	Dhanush Lin	lingewo@usc.edu	Yash Thakur	ythakur@usc.edu
7	NLP Anomaly Detection Library	jiqil(Bourne) Li	ji77629@usc.edu	Yuangang Li	yuangang@usc.edu	Leo Lee	ltee@usc.edu	Bowen Wang	bwang443@usc.edu	Huixian Gong	huixian@usc.edu
8	Medical Advice Chatbot	Dihan Li	dihani@usc.edu	Hantao Ye	hantao@usc.edu	Kai Zheng	kzheng44@usc.edu	Zhenyi Xie	zhenyi@usc.edu		
9	Untitled	Antonio Revilla	honggui@usc.edu	Govind Thakur	gthakur@usc.edu	Ishita Mehta	ishitame@usc.edu	Sachin Kumar	skumar42@usc.edu	Venkata Sai Suryasudha	vsudu@usc.edu
10	Automated PII Detection and Removal	Gokhan Mungan	mungan@usc.edu	Jonathan Kaspi	kaspiris@usc.edu	Tianyi Zhao	tzhao566@usc.edu	Neeraj Iyer	n.neeraj@usc.edu	Aaron Dardik	dardik@usc.edu
11	LLMs for Non-Generative Code Clone Detection Task	Ananya Kotha	akotha@usc.edu	Apoorva Sharma	asharma2@usc.edu	Saloni Oswal	smoswal@usc.edu	Gurunadh Pamarti	pamarti@usc.edu		
12	Cautious Reasoner LM for Coding Tasks	Clement Chan	cjchan@usc.edu	Bruno Segovia	segovia@usc.edu	Scott Susanto	scottsus@usc.edu	Trevor Asber	asbery@usc.edu		
13	GPTE self-augmenting fine-tuning	Ziyan Zeng	ziyanzen@usc.edu	Vincent Deng	vdeng@usc.edu	Junyang Cai	cajunya@usc.edu	hanjin@usc.edu		Han Xiao	hxiao603@usc.edu
14	Synthetic Satellite Image Generation for Forecasting Urban and Natural Changes	Jay Prajapati	jprajapa@usc.edu	Adwait Thergaonkar	thergaon@usc.edu	Cynthia Chang	cchang80@usc.edu	Chanavi Siricharanavis	chanavis@usc.edu	Shourya Kothari	sakothar@usc.edu
15	Using Deep Learning for Image Compression	Mahika Moshuni	moshuni@usc.edu	Diana Pham	dianaph@usc.edu	Karen Wang	kwang760@usc.edu	Praise Olukile	olukile@usc.edu	Vedant Modi	vkmodi@usc.edu
16	Analyzing Community Opinion Dynamics in Response to Significant Sociopolitical Events with Large Language Models	Kshitij Pawar	kshitijv@usc.edu	Syed's Masoon	syedmasoon@usc.edu	Sanjay Raghavend	sraghave@usc.edu	Mihika Gaonkar	mihikarp@usc.edu		
17	Investigating the Impact of Differing Word Embeddings on Gender Bias in Language Models (LLMs)	Omi Wakode	wakode@usc.edu	Padmaja Borwankar	borwanka@usc.edu	Yatharth Arora	yarora@usc.edu	Andrew Klotz	aklotz@usc.edu	Trathi Rathi	trathi@usc.edu
18	Self Contradictory reasoning and logical fallacies detection	Ritesh Reddy Bandla	rbandam@usc.edu	Paritosh Kadar	pskadam@usc.edu	Sandeep Menon	menons@usc.edu	Aviraj Goel	aviralgo@usc.edu	Krishna Wankhede	wankhede@usc.edu
19	Benchmarking unsupervised environment design (UED) algorithms	Nikhil Sathyendaray	nsathyan@usc.edu	Prajna Narayan Hollakun	holakun@usc.edu	Ronak Jain	ronakjai@usc.edu	Akshay Kaushik	akshay@usc.edu	Rithvik Mohan	vi.vijalapu@usc.edu
20	Untitled	Rishabh Shinde	shinder@usc.edu	Pallavi Udatewar	udatewar@usc.edu	Aryan Gandhi	aryangan@usc.edu	Aditya Jadhav	adityadh@usc.edu		

# Project Check-in

Each project has 2 mandatory check-in with your TA

- Potentially one before the mid-report due (March 22nd)
- Potentially another one before final report is due (May 3rd)
- We will work out the logistic then

We have 55 projects – each TA is in charge of 7-8 projects

- They will be your point of contact for project questions
- Use Piazza, office hours, or scheduled appointment with them to make sure your questions are addressed

# Cloud Credits

First, **Google** has \$300 for all new signups --

<https://cloud.google.com/free?hl=en>

- Also, we have applied for the Google Cloud credit (\$50/student) -- considering pooling them for the project.

I have also applied for **USC advanced computing credit** --

<https://www.carc.usc.edu/> to request the usage of it, please add your information

- <https://docs.google.com/spreadsheets/d/1KWPIHyMVZ-2o47IhpFC9gxDYusbchPuvX971HhtUULI/edit?usp=sharing>

# Cloud Credits

## Google Cloud Credit:

<https://vector.my.salesforce-sites.com/GCPEDU?cid=U4cYmFRYRNMCY3QWOi3mu7niNVKxs%2Fq76pJeSiD9VVD23kqHSaeE71KE2OJOYly6/>

**USC Credit:** <https://docs.google.com/spreadsheets/d/1KWPIHyMVZ-2o47IhpFC9gxDYusbchPuvX971HhtUULI/edit#gid=0>

See details: <https://piazza.com/class/loendg1jxkryr/post/21>

# Assignments

## Designed by TAs:

- neural network, decision tree, and random forests
- Note Python is a prerequisite of the course
- PyTorch has a 60 min quick tutorial:  
[https://pytorch.org/tutorials/beginner/deep\\_learning\\_60min\\_blitz.html](https://pytorch.org/tutorials/beginner/deep_learning_60min_blitz.html)

We leverage Gradescope similarity check + point check:

- Also potentially my research, which is?
- Do the homework independently while general discussion is allowed
- ChatGPT is not recommended, but hard to check though...

# Assignments

## Due in 3 weeks...

- The project should still be ongoing
- Questions on piazza will be addressed by TA
  - Yuehan, Varun, and Ziyi will help
- Also TA hours will be helpful for all the Tas

Do not do it in the last minute – recall we do not do late days due to the size of the course

# Interaction on Piazza

**Please post general questions to all students:**

- You will get your answer much faster!
- Again – do not post your assignment answer, but the general questions are fine!

If you actively answer others' questions, we will take note of your name  
- if you need a bump by the end of the semester, we will get you there  
automatically due to your contribution to the course!

# TA Office Hours

We have a slightly easier way to handle office hours (not that ideal)

 i **Zihao He** 4 days ago Actions ▾

Hi everyone, we will share our TA OHs in here in the future.

	Varun Bhatt	Zihao He	Zizhao Hu	Ayush Jain	Ziyi Liu	Yuehan Qin	Pengda Xiang
Week of Feb 5	Tuesday 14:00 - 15:00, RTH 419 (preferred) or <a href="#">Zoom</a> ,	Thu, 2-3pm Leavey 202C, <a href="#">Zoom Link</a>	Wed, 10-11am Leavey, 201F <a href="#">Zoom</a>	Wed, 2-3pm RTH 4th floor lobby		Tuesday 10:00-11:00 <a href="#">Zoom</a>	
Week of Feb 12							

[good comment](#) | 2

# Midterm – Note there is no Final!

Feb 23<sup>rd</sup> – **in-class – paper-based** – open books (although I have no idea which books to use)

- 50 + 6 (bonus) multiple choice questions + single answer questions, e.g., a single number
- there is no make-up unless you have medical emergency with proof ☺
- 1 bonus for each quiz and three bonus for the exam (0.5 each so it is 6 questions)

# Stop Talking during the Lecture

question @58 ⚡ ★ 🔒

Copy @58

## Stop talking during class when the professor is speaking

It is extremely annoying when people are talking when the professor is speaking. Half the time I can barely make out the words that the professor going to talk do it outside or don't come to class.

logistics

Edit good question | 45 Updated 6 days ago by Michael Kingsley (Ar)

S the students' answer, where students collectively construct a single answer +1

Edit good answer | 1 Updated 6 days ago by Chanavi Singh (Ano)

i the instructors' answer, where instructors collectively construct a single answer



- Ask these students to share their questions/topics on the stage
- Also, I will start walking around during the lecture time to help address these questions more closely (for the benefit of my neck)

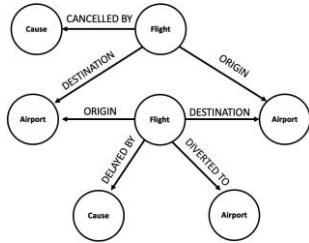
# Graph Neural Networks (GNN)

## Basics

# Why Graphs?

Graphs are a general language for describing and analyzing entities with relations/interactions

# Many Types of Data are Graphs (1)

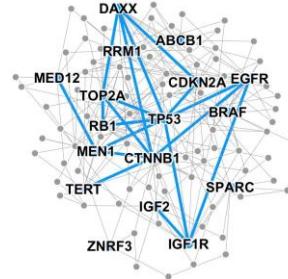


**Event Graphs**



Image credit: [SalientNetworks](#)

**Computer Networks**



**Disease Pathways**

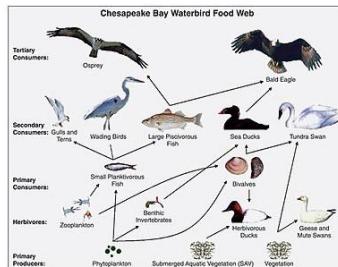


Image credit: [Wikipedia](#)

**Food Webs**



Image credit: [Pinterest](#)

**Particle Networks**



Image credit: [visitlondon.com](#)

**Underground Networks**

# Many Types of Data are Graphs (2)



Image credit: [Medium](#)

**Social Networks**

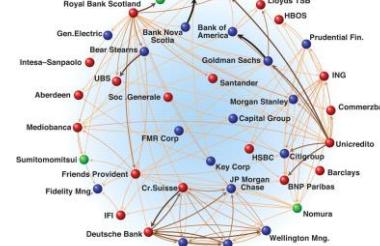


Image credit: [Science](#)

**Economic Networks**

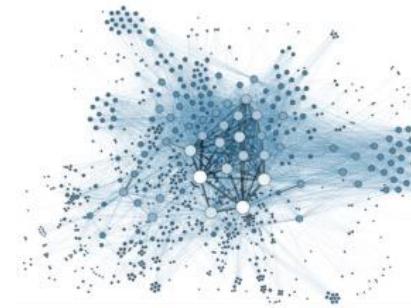
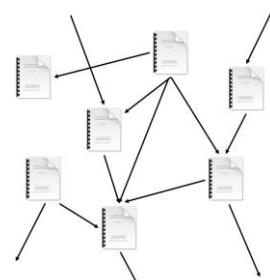


Image credit: [Lumen Learning](#)

**Communication Networks**



**Citation Networks**



Image credit: [Missoula Current News](#)

**Internet**

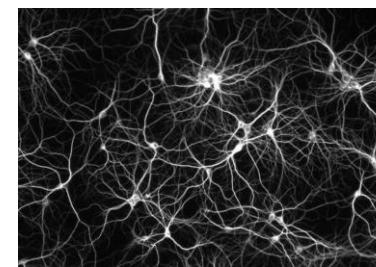


Image credit: [The Conversation](#)

**Networks of Neurons**

# Many Types of Data are Graphs (3)

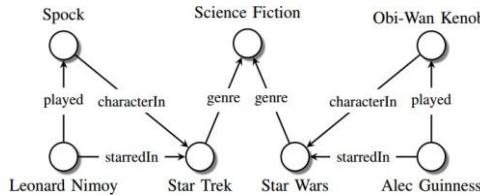


Image credit: [Maximilian Nickel et al](#)

## Knowledge Graphs

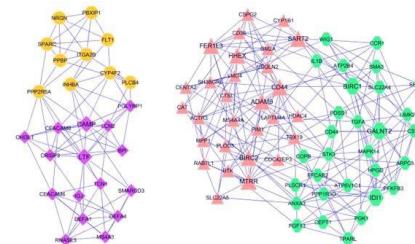


Image credit: [ese.wustl.edu](#)

## Regulatory Networks

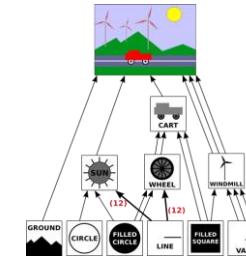


Image credit: [math.hws.edu](#)

## Scene Graphs

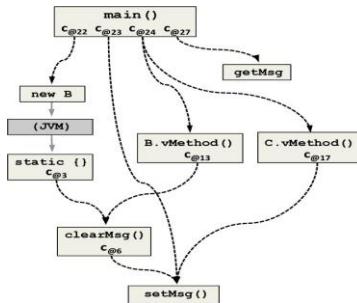


Image credit: [ResearchGate](#)

## Code Graphs

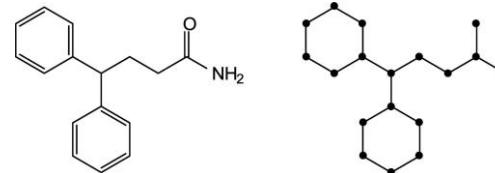


Image credit: [MDPI](#)

## Molecules

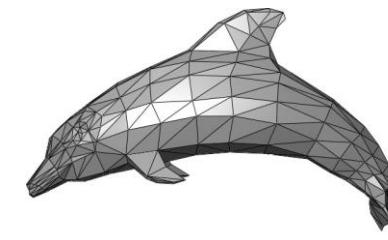


Image credit: [Wikipedia](#)

## 3D Shapes

# Graphs: Machine Learning

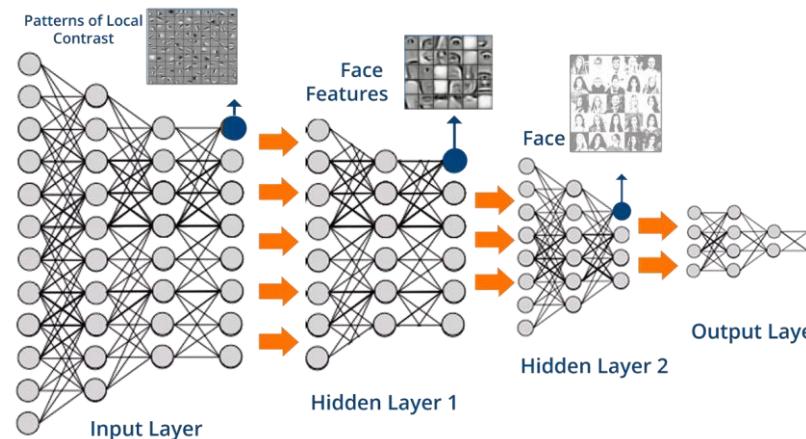
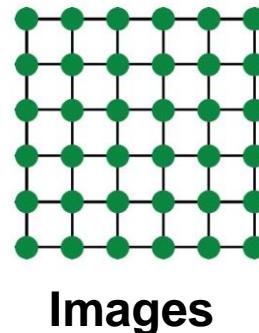
Complex domains have a rich relational structure, which can be represented as a **relational graph**

By explicitly modeling  
relationships, we  
achieve better performance!

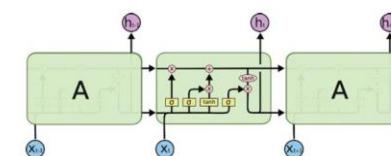
Main question:  
How do we take advantage of relational  
structure for better prediction?

# Today: Modern ML Toolbox

CNN



Text/Speech



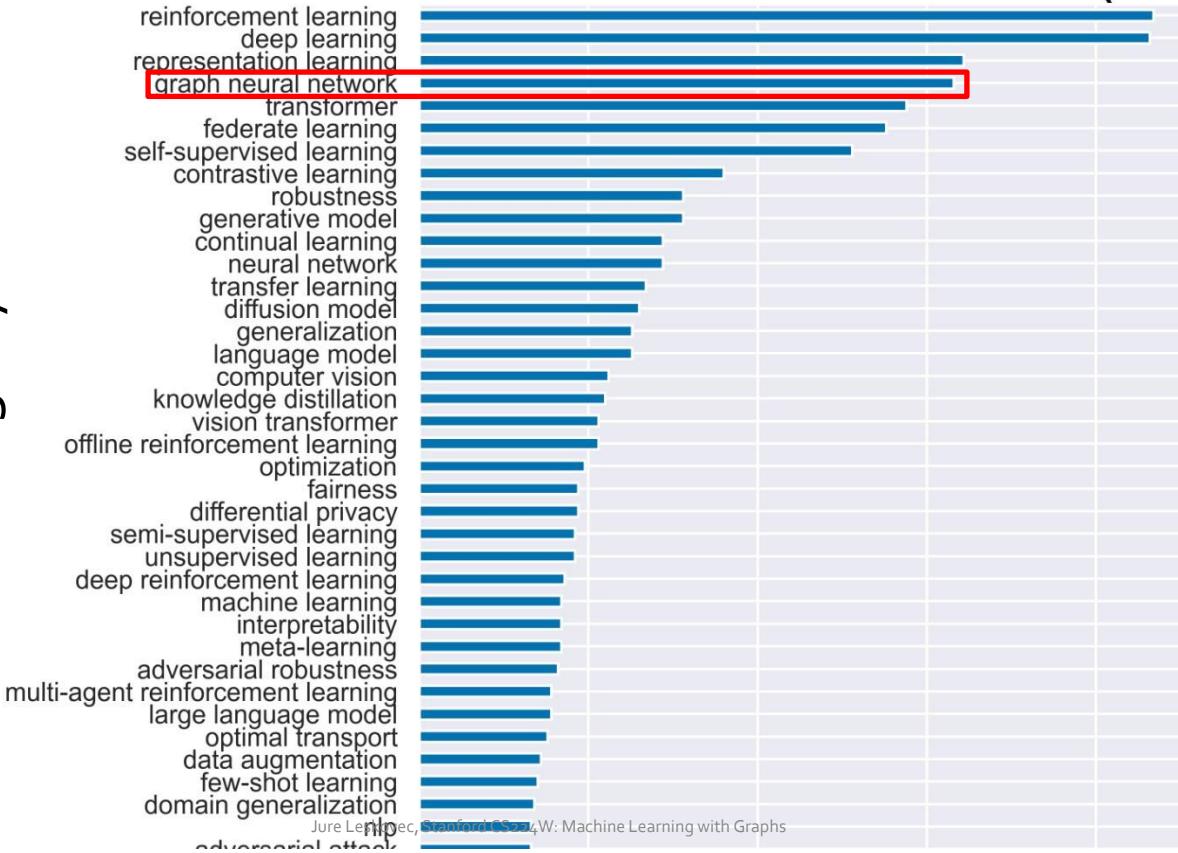
RNN/  
LSTM

Modern deep learning toolbox is designed  
for simple sequences & grids

# Hot subfield in ML

## 50 MOST APPEARED KEYWORDS (2023)

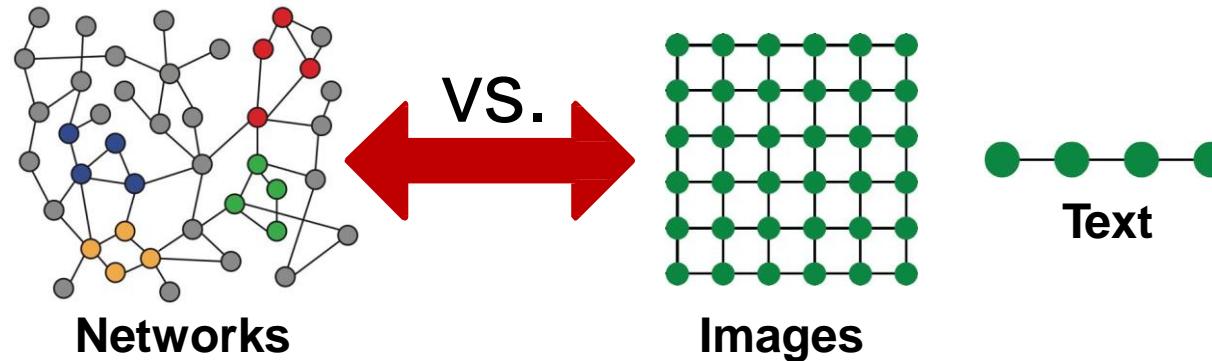
ICLR 2023 keywords



# Why is Graph Deep Learning Hard?

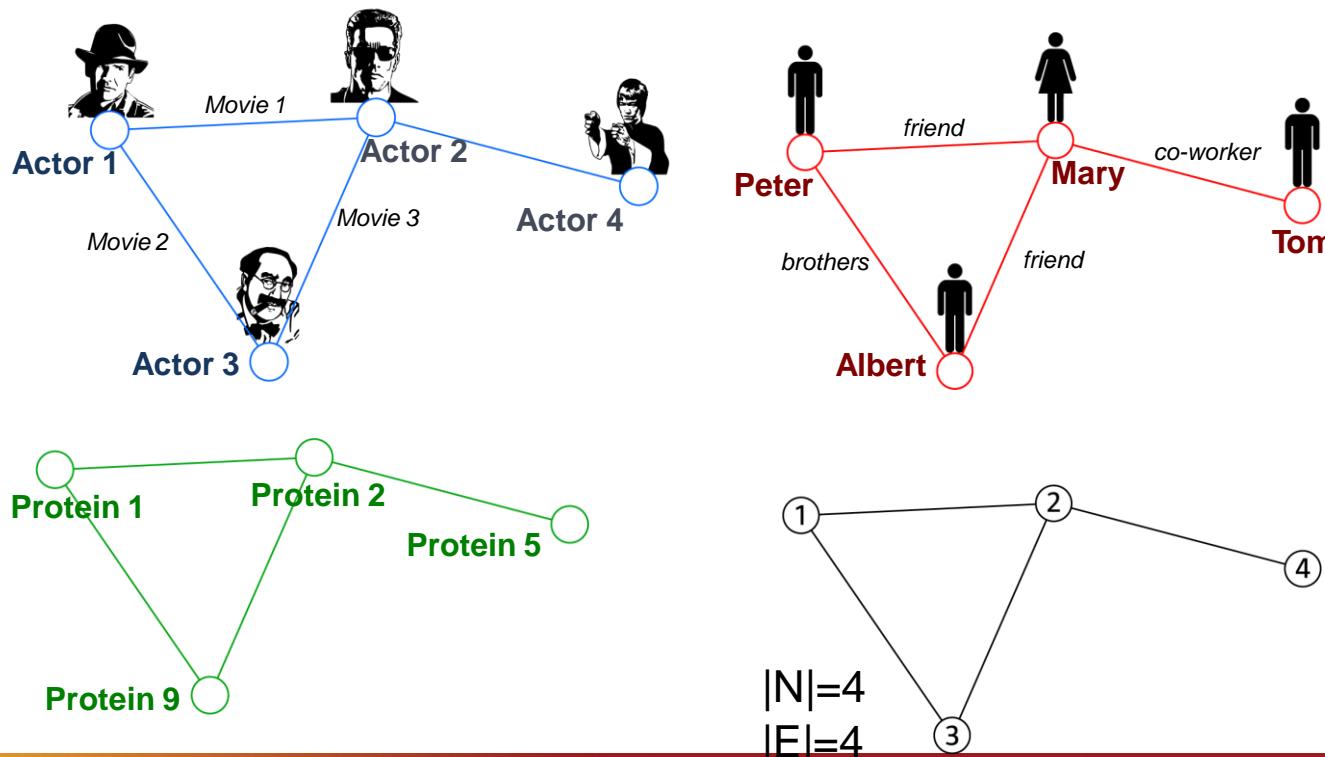
## Networks are complex:

- Arbitrary size and complex topological structure (*i.e.*, no spatial locality like grids)



- No fixed node ordering or reference point
- Often dynamic and have multimodal features

# Graphs: A Common Language



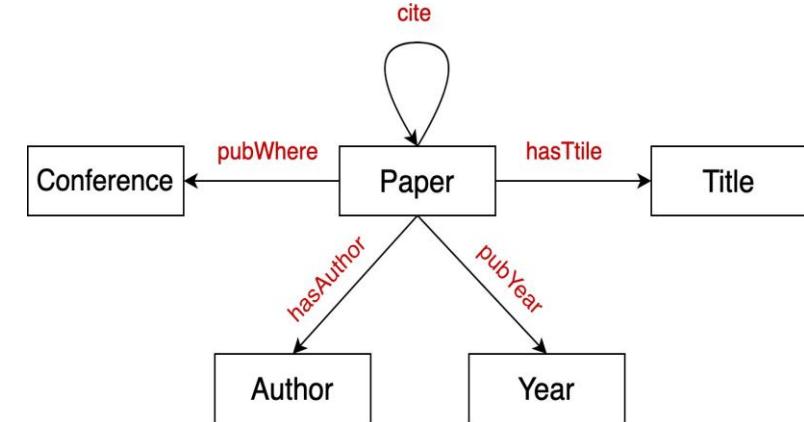
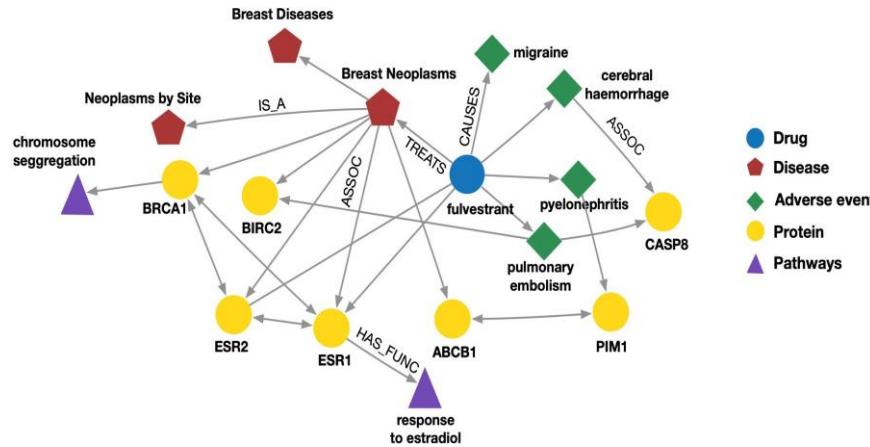
# Heterogeneous Graphs

- A heterogeneous graph is defined as

$$G = (V, E, R, T)$$

- Nodes with node types  $v_! \in V$
- Edges with relation types  $(v_!, r, v^") \in E$
- Node type  $T(v_!)$
- Relation type  $r \in R$
- Nodes and edges have attributes/features

# Many Graphs are Heterogeneous



## 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

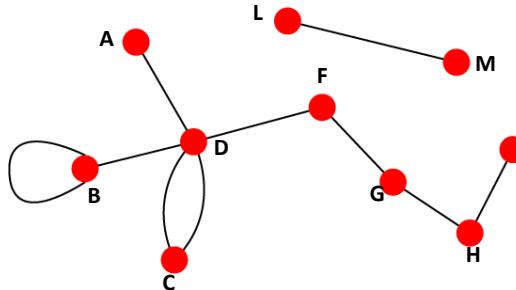
# Choosing a Proper Representation

- **How to build a graph:**
  - What are nodes?
  - What are edges?
- **Choice of the proper network representation of a given domain/problem determines our ability to use networks successfully:**
  - In some cases, there is a unique, unambiguous representation
  - In other cases, the representation is by no means unique
  - The way you assign links will determine the nature of the question you can study

# Directed vs. Undirected Graphs

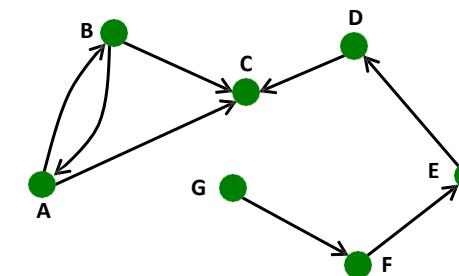
## Undirected

- Links: undirected  
(symmetrical, reciprocal)



## Directed

- Links: directed

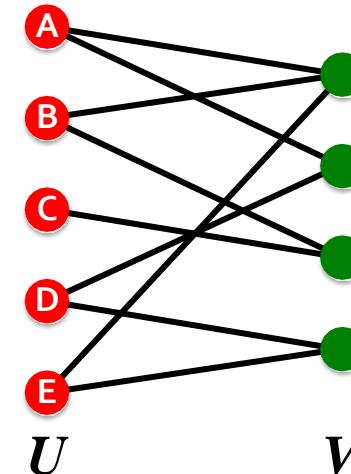


## Other considerations:

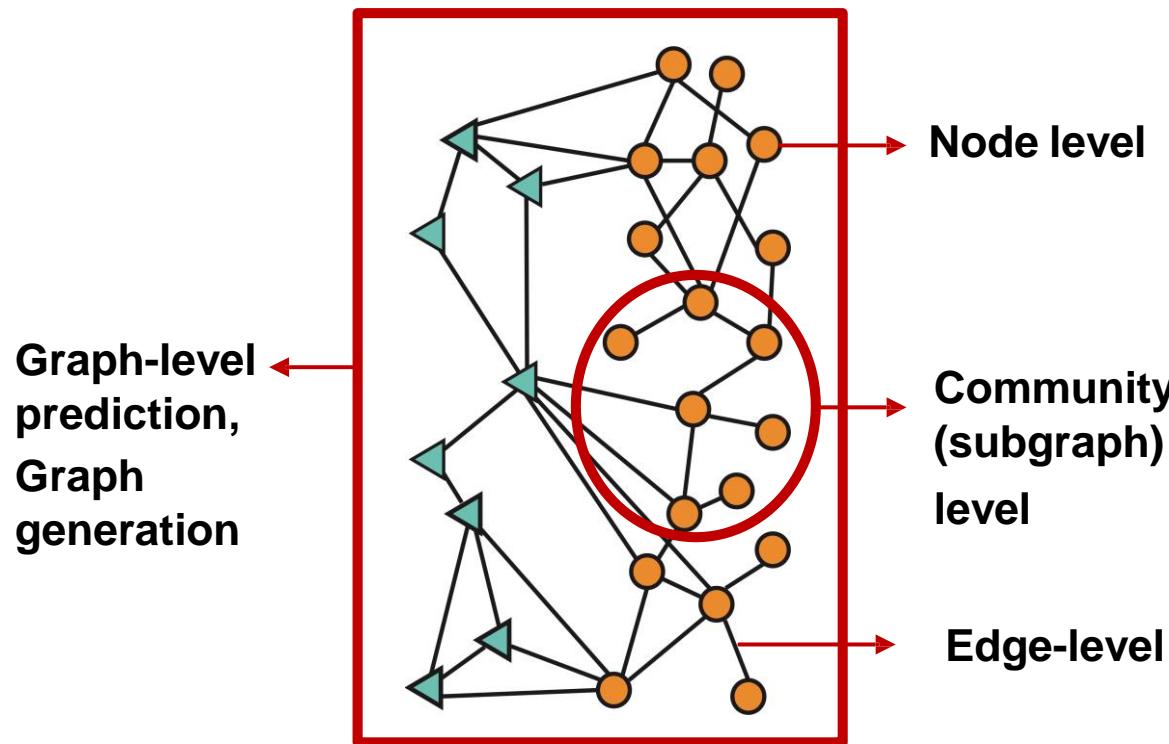
- Weights
- Properties
- Types
- Attributes

# Bipartite Graph

- **Bipartite graph** is a graph whose nodes can be divided into two disjoint sets  $U$  and  $V$  such that every link connects a node in  $U$  to one in  $V$ ; that is,  $U$  and  $V$  are **independent sets**
- **Examples:**
  - Authors-to-Papers (they authored)
  - Actors-to-Movies (they appeared in)
  - Users-to-Movies (they rated)
  - Recipes-to-Ingredients (they contain)

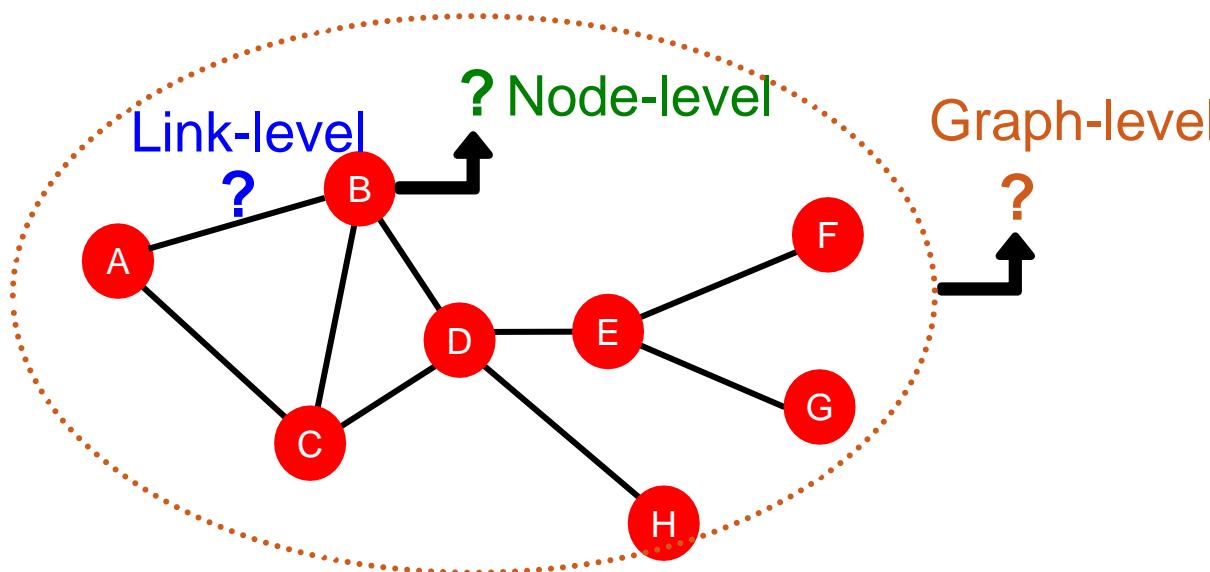


# Different Types of Tasks

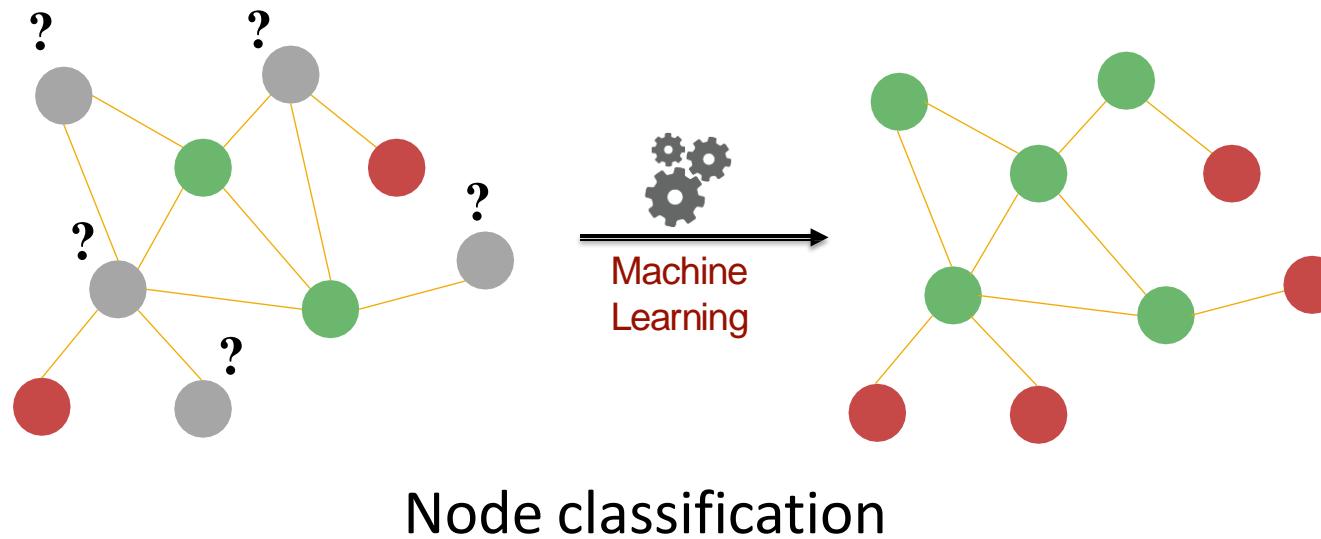


# Machine Learning Tasks: Review

- Node-level prediction
- Link-level prediction
- Graph-level prediction



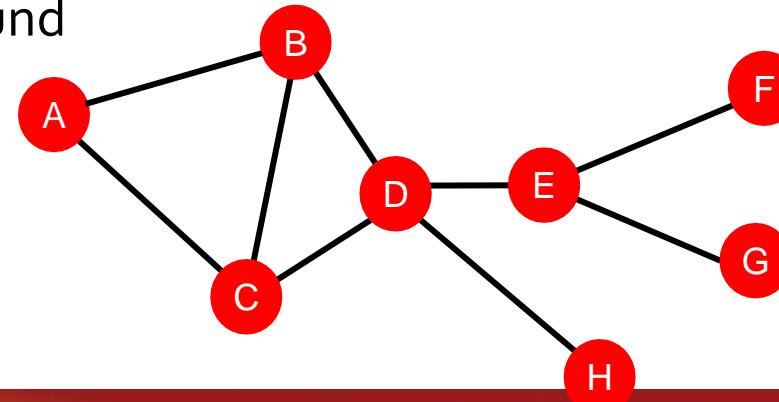
# Node-Level Tasks



# Node-Level Network Structure

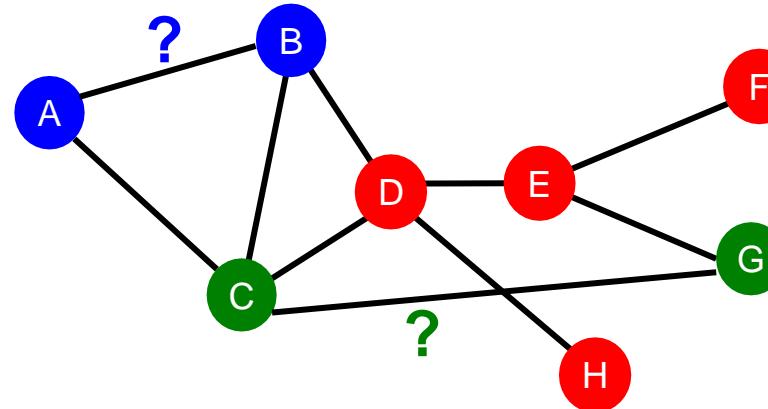
**Goal:** Characterize the structure and position of a node in the network:

- Node degree
- Node importance & position
  - E.g., Number of shortest paths passing through a node
  - E.g., Avg. shortest path length to other nodes
- Substructures around the node



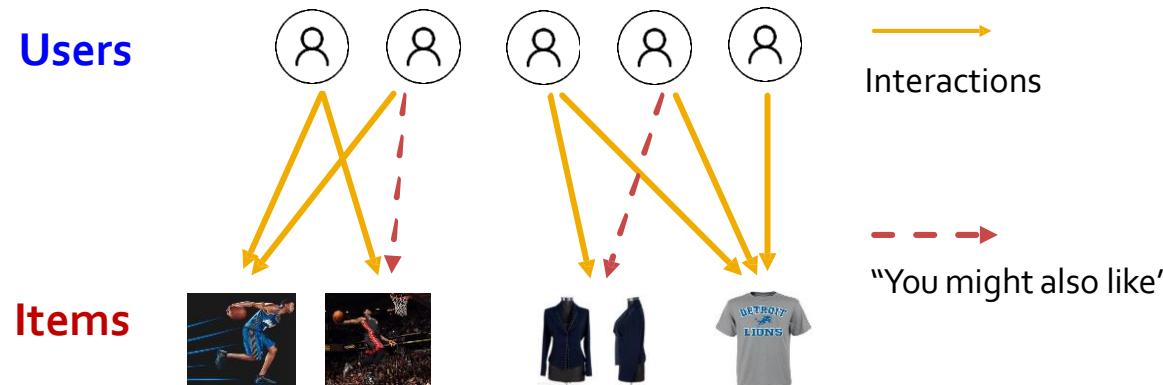
# Link-Level Prediction Task

- The task is to predict **new/missing/unknown links** based on the existing links.
- At test time, node pairs (with no existing links) are ranked, and top  $K$  node pairs are predicted.
- Task: Make a prediction for a pair of nodes.



# Example (1): Recommender Systems

- **Users interacts with items**
  - Watch movies, buy merchandise, listen to music
  - **Nodes:** Users and items
  - **Edges:** User-item interactions
- **Goal: Recommend items users might like**



# PinSage: Graph-based Recommender

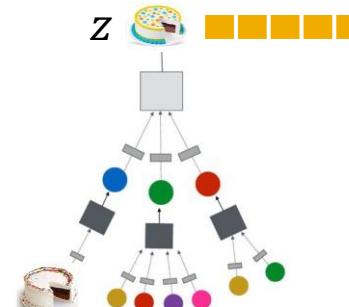
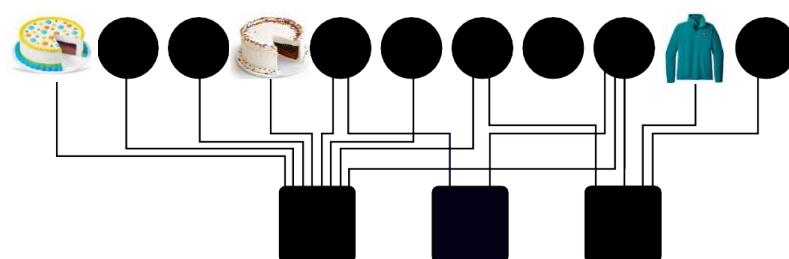
**Task:** Recommend related pins to users



**Task:** Learn node embeddings  $z_i$  such that

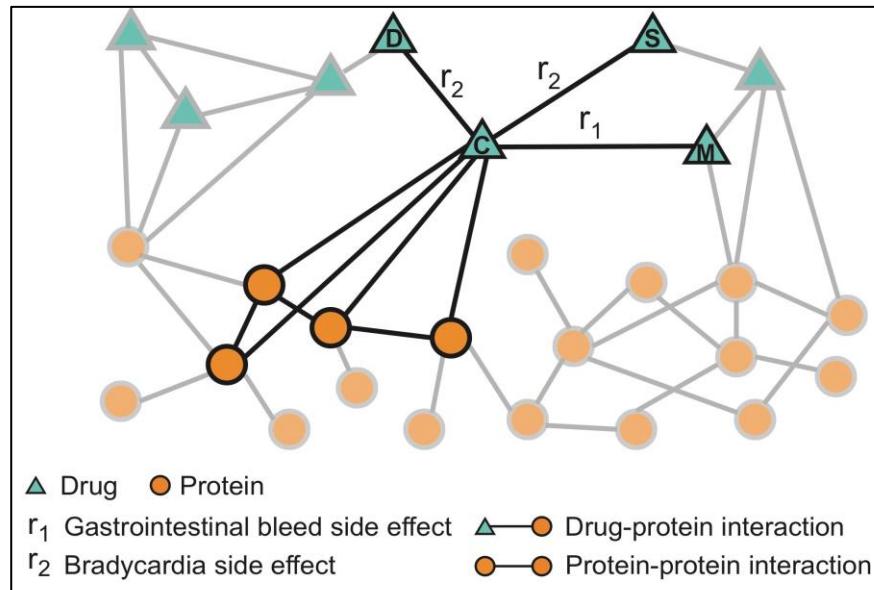
$$d(z_{cake1}, z_{cake2}) < d(z_{cake1}, z_{sweater})$$

Predict whether two nodes in a graph are related

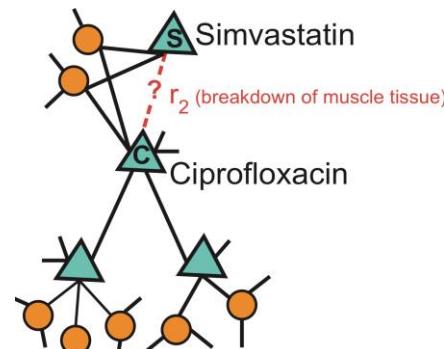


# Biomedical Graph Link Prediction

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

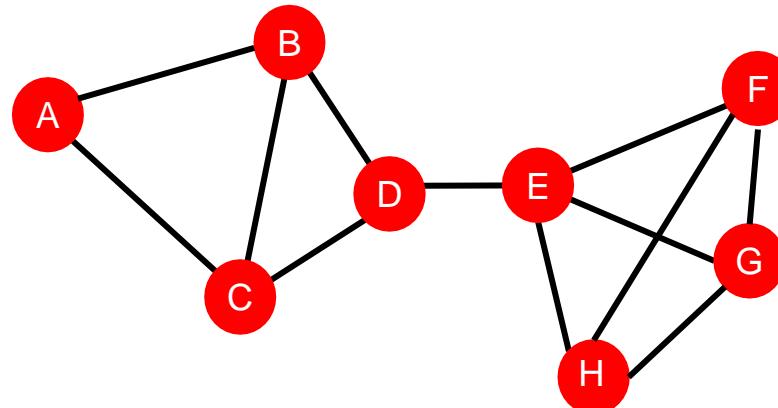


**Query:** How likely will Simvastatin and Ciprofloxacin, when taken together, break down muscle tissue?



# Graph-Level Features

- **Goal:** We want make a prediction for an entire graph or a subgraph of the graph.
- **For example:**



# Example (1): Traffic Prediction

**Best** 17 min | 1 hr 46 | 3 hr 56 | 1 hr 5

**Gas** **EV charging** **Hotels** **More**

University of California, Los Angeles  
USC Viterbi School of Engineering, 3650 | Add destination

Leave now | Options

Send directions to your phone

**via I-10 E** **17 min** | 14.0 miles

Fastest route now due to traffic conditions

Details

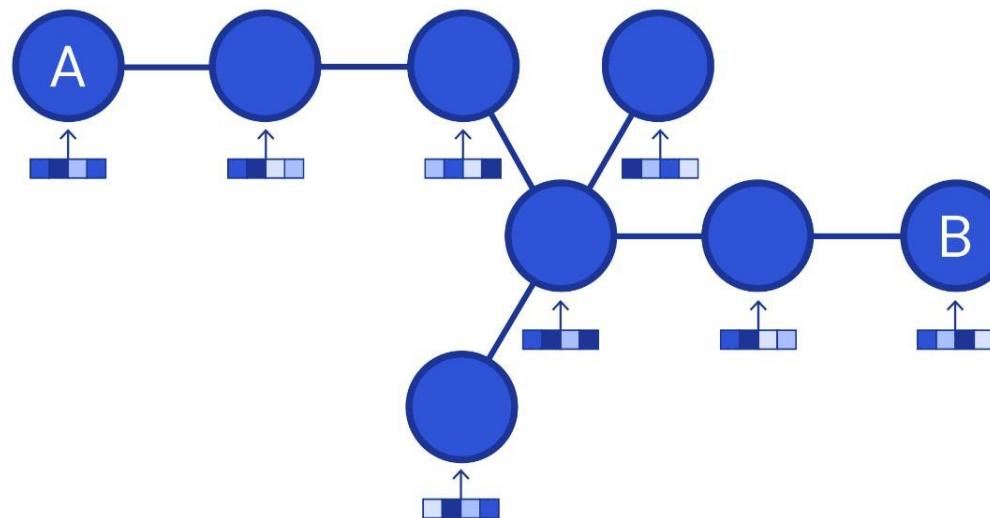
Explore nearby USC Viterbi School of Engineering

Restaurants | Hotels | Gas stations | Parking Lots | More

Map data ©2024 United States Terms Privacy Send Product Feedback 1 mi 11:21 PM

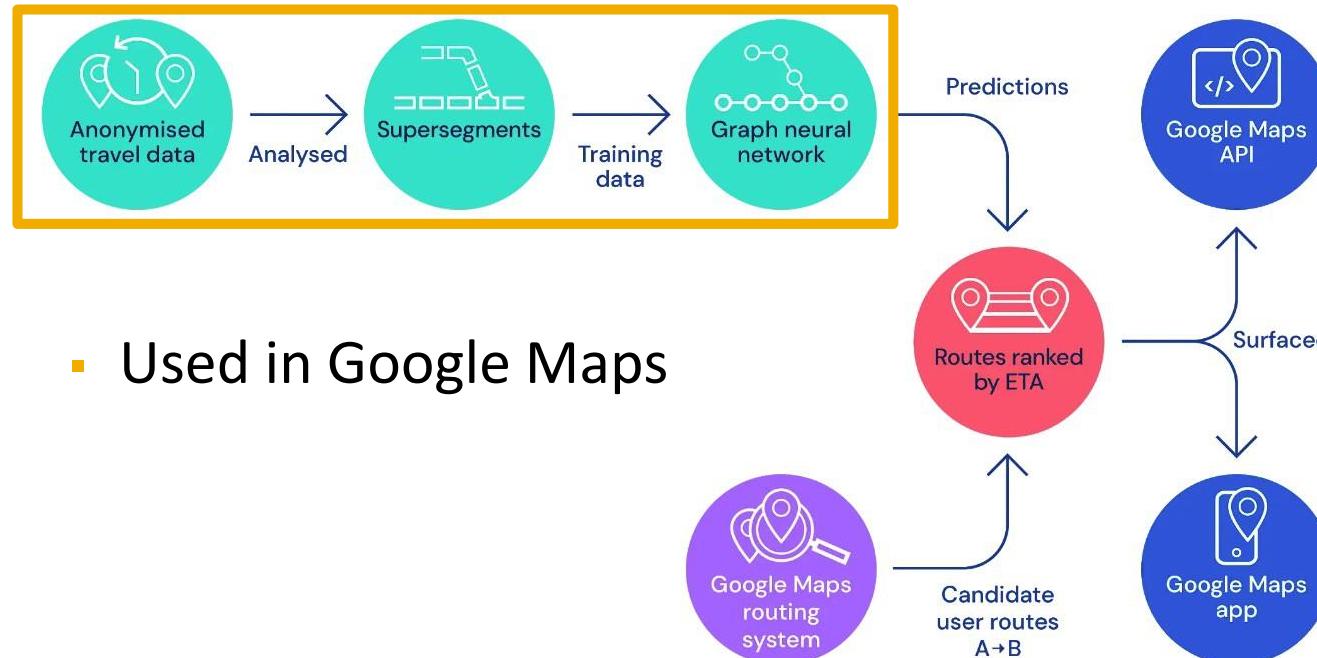
# Road Network as a Graph

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



# Traffic Prediction via GNN

## Predicting Time of Arrival with Graph Neural Networks

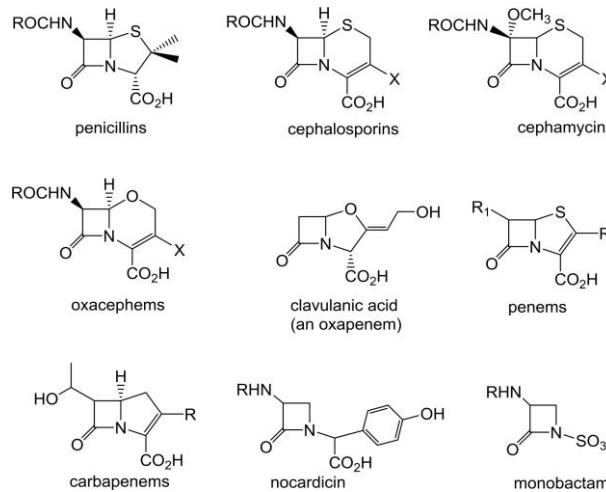


THE MODEL ARCHITECTURE FOR DETERMINING OPTIMAL ROUTES AND THEIR TRAVEL TIME.

Image credit: [DeepMind](#)

# Example (2): Drug Discovery

- Antibiotics are small molecular graphs
  - **Nodes:** Atoms
  - **Edges:** Chemical bonds

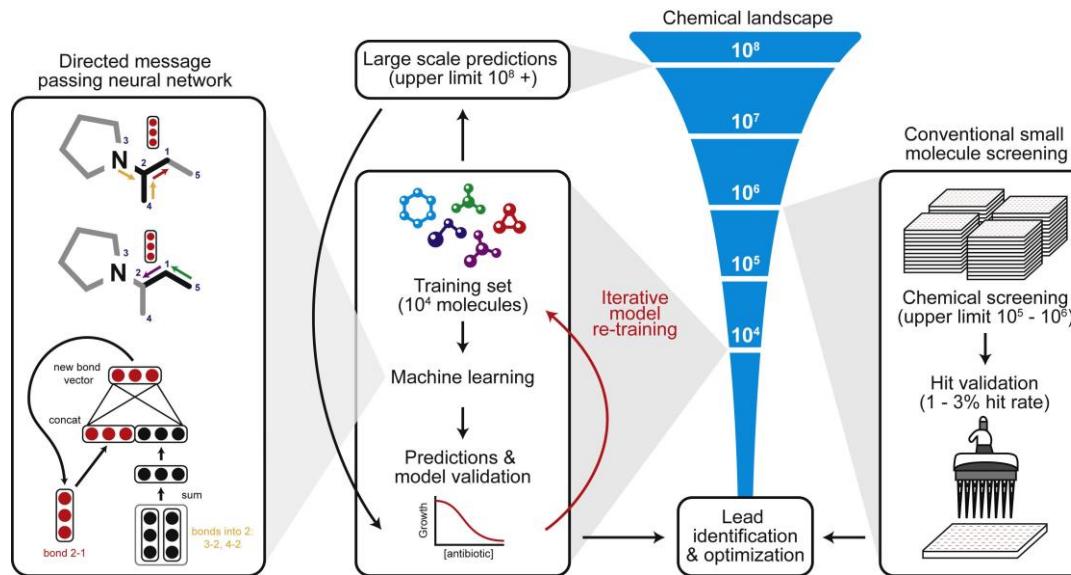


Konaklieva, Monika I. "Molecular targets of  $\beta$ -lactam-based antimicrobials: beyond the usual suspects." *Antibiotics* 3.2 (2014): 128-142.

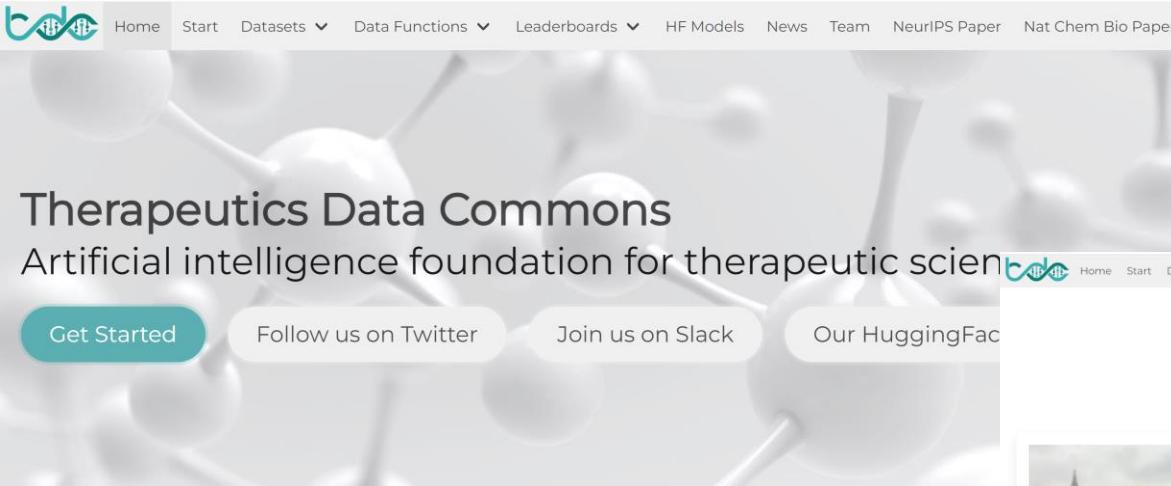
Image credit: [CNN](#)

# Deep Learning for Antibiotic Discovery

- A Graph Neural Network **graph classification model**
- Predict promising molecules from a pool of candidates



Stokes, Jonathan M., et al. "A deep learning approach to antibiotic discovery." *Cell* 180.4 (2020): 688-702.



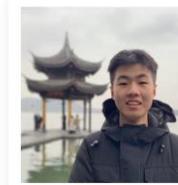
# Therapeutics Data Commons

Artificial intelligence foundation for therapeutic science

[Get Started](#)[Follow us on Twitter](#)[Join us on Slack](#)[Our HuggingFace](#)

## Core Team

Developing and Maintaining Therapeutics Data Commons

[Contact us](#)

Kexin Huang  
Stanford  
[Website](#)



Tianfan Fu  
Georgia Tech  
[Website](#)



Wenhao Gao  
MIT  
[Website](#)



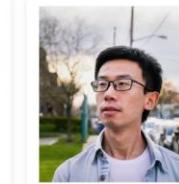
Yuanqi Du  
Cornell  
[Website](#)



Ada Fang  
Harvard  
[Website](#)



George Dasoulas  
Harvard  
[Website](#)



Yue Zhao  
CMU  
[Website](#)



Nitin Pasumarty  
LinkedIn  
[Website](#)

Spring' 24 Y. Zhao

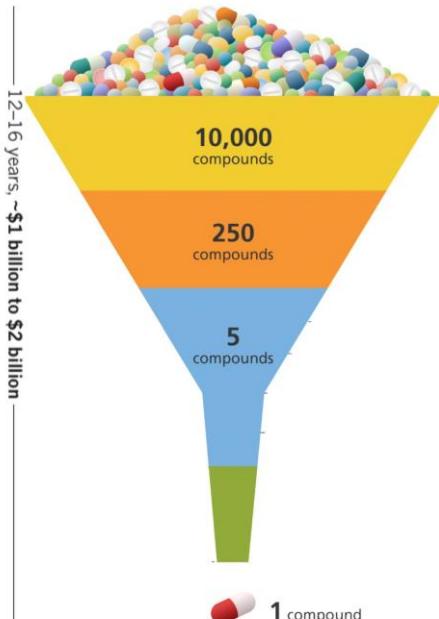
Viterbi  
School of Engineering

# Challenges in Drug Discovery & Development

High cost, long time

Various and emerging diseases

Abundant molecular modality



Step 1: Target Identification



Step 2: Hit Identification and Lead Optimization



Step 3: Clinical Research



Step 4: FDA Review



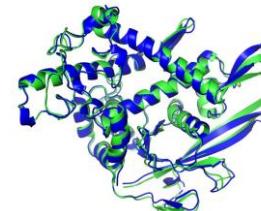
Step 5: Manufacturing and Post-marketing



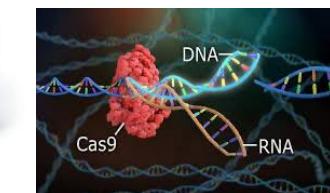
Small molecules



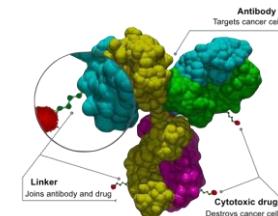
Vaccines



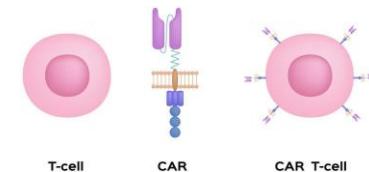
Proteins



Gene-editing



Anti- Nano-bodies



T-cell

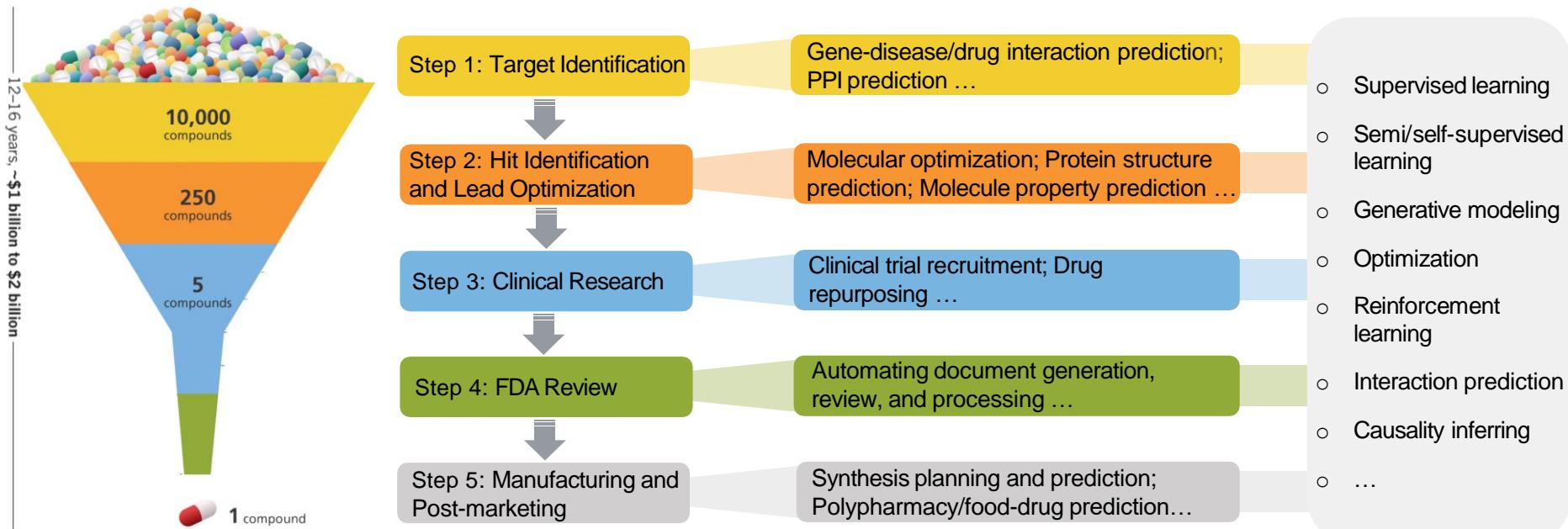
CAR

CAR-T cell

CAR-T

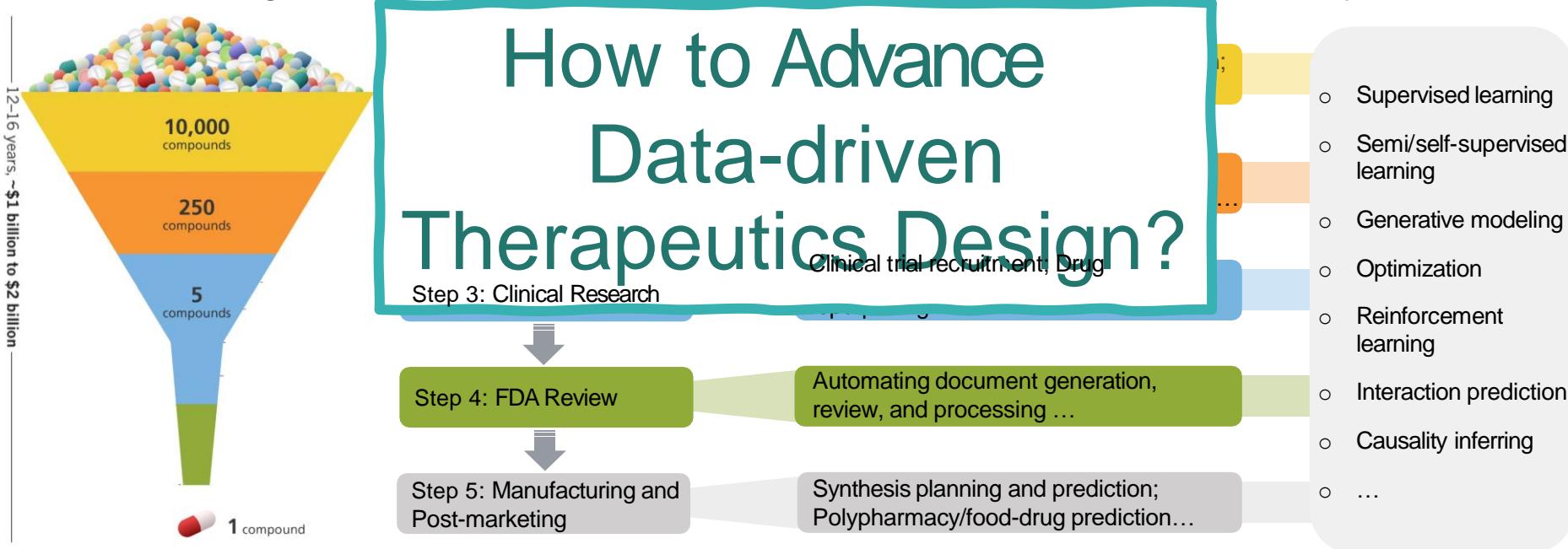
# Challenges of AI in Drug Discovery & Development

- Formulate various tasks as machine-learning-solvable tasks.
- Identify, retrieve, and process datasets of many different types scattered around.
- Assess algorithmic advances to align with real-world and clinical deployment.



# Challenges of AI in Drug Discovery & Development

- Formulate various tasks as machine-learning-solvable tasks.
- Identify, retrieve, and process datasets of many different types scattered around.
- Assess algorithmic advances to align with real-world and clinical deployment.



# Therapeutics Data Commons



THERAPEUTICS  
DATA COMMONS

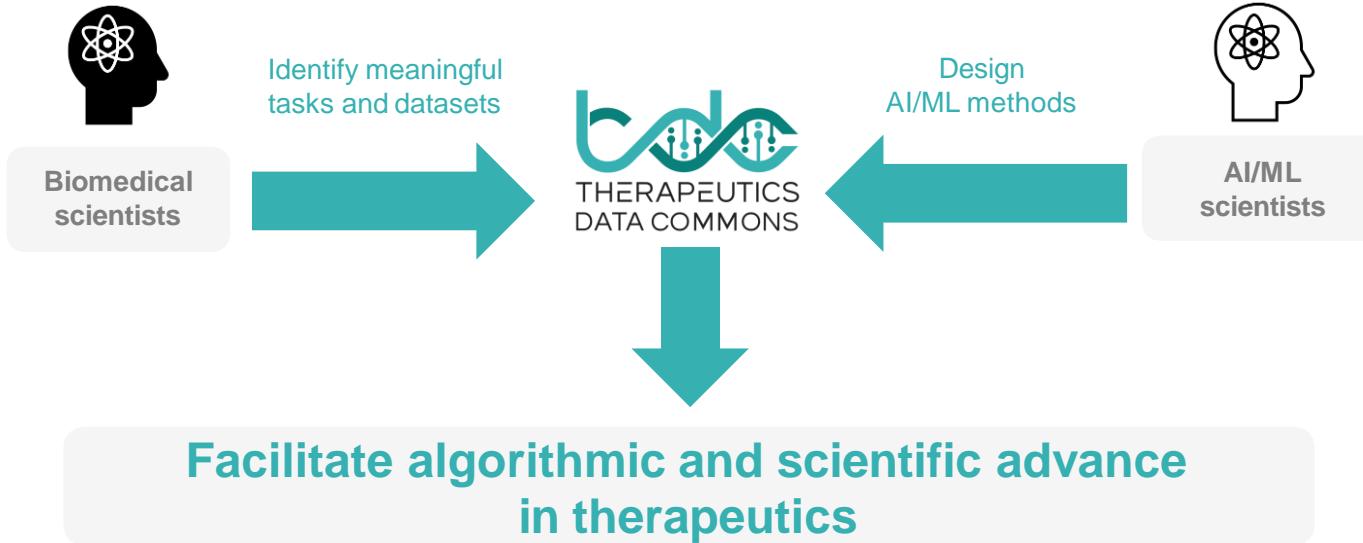
Website: <https://tdcommons.ai> (or QR code)



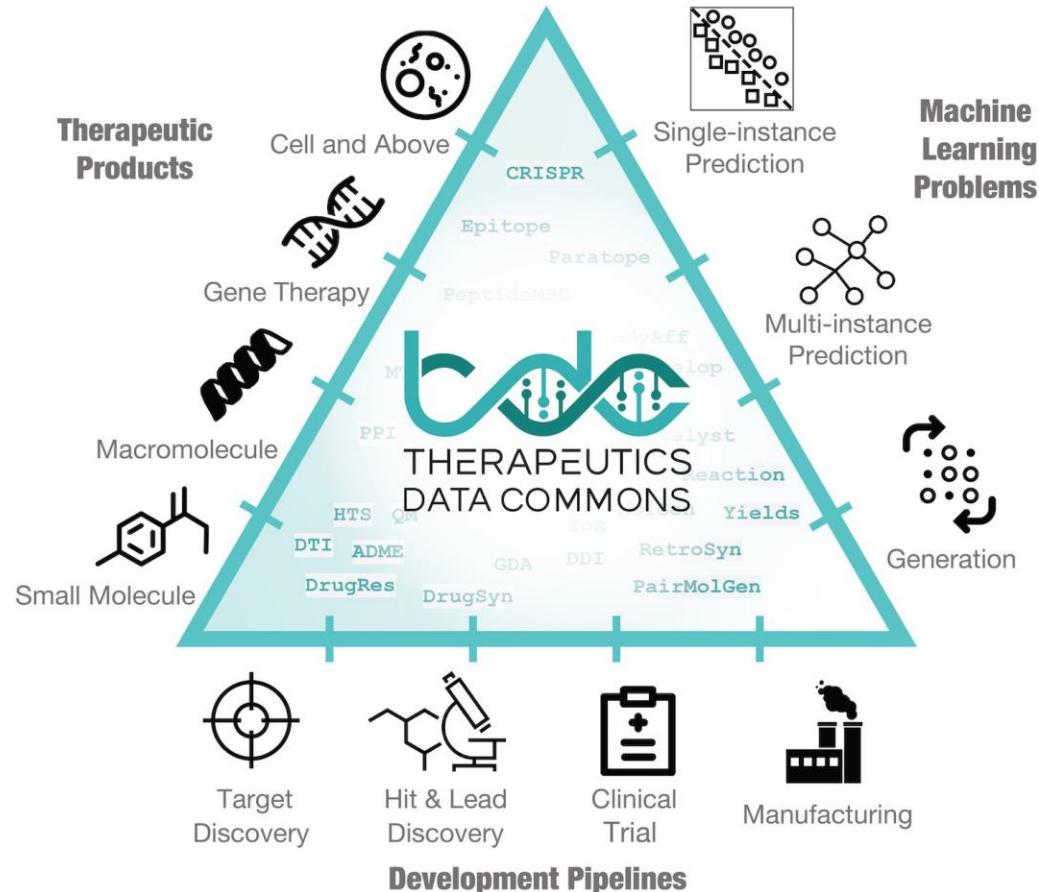
Paper: <https://www.nature.com/articles/s41589-022-01131-2>



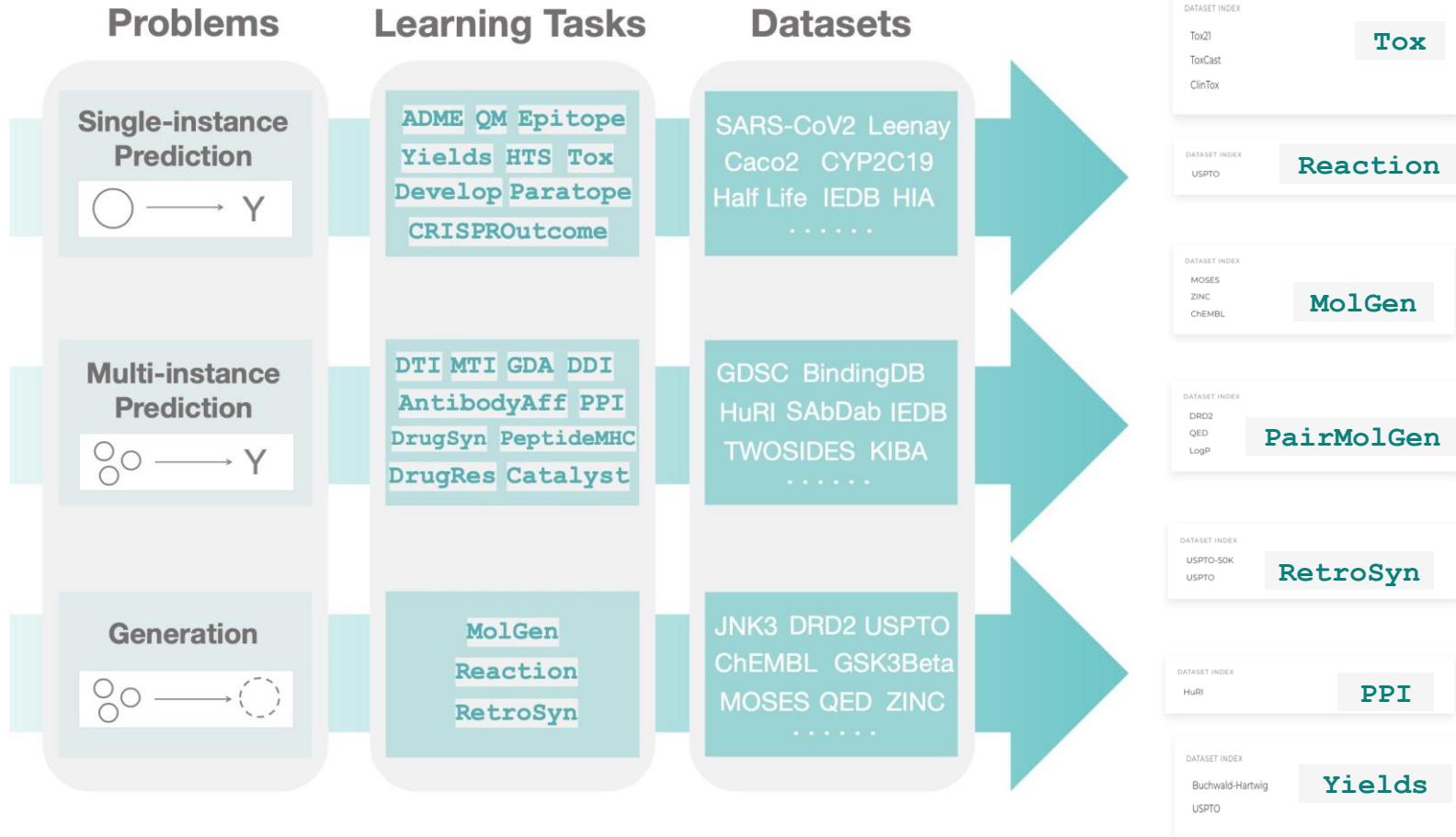
Github: <https://github.com/mims-harvard/TDC>



# Wide Range of Therapeutic Modalities and Pipeline



# Three-Tier Design

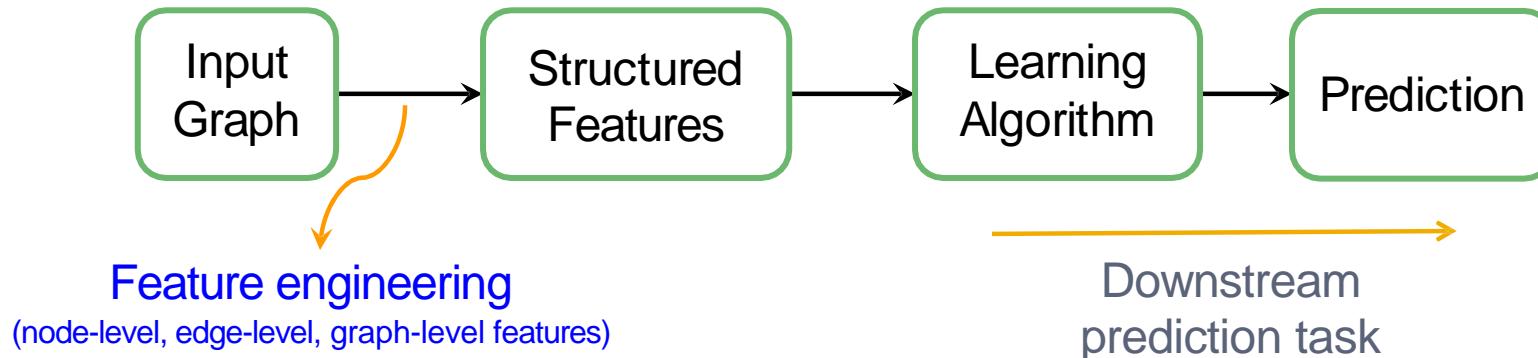


# **Graph Neural Networks (GNN)**

**Details**

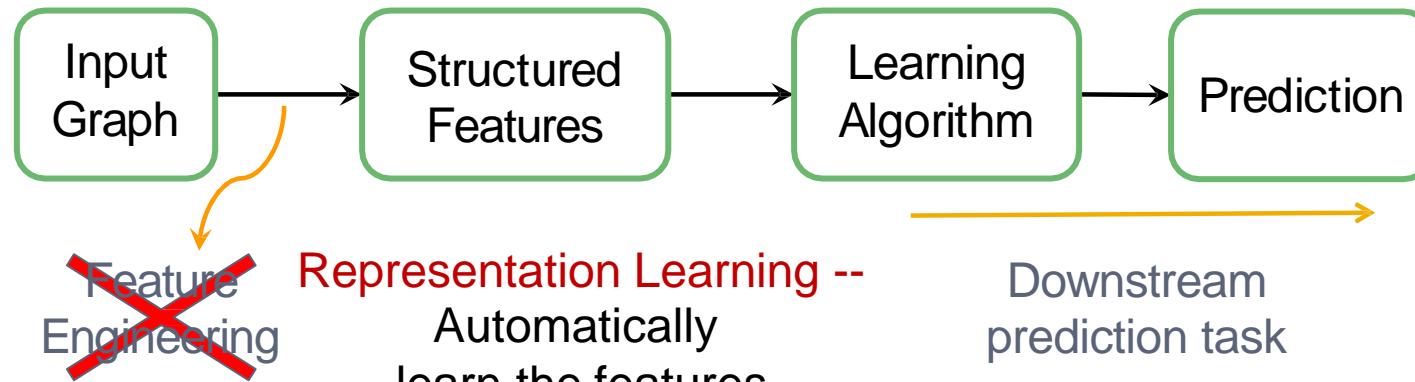
# Recap: Traditional ML for Graphs

Given an input graph, extract node, link and graph-level features, then learn a model (SVM, neural network, etc.) that maps features to labels.



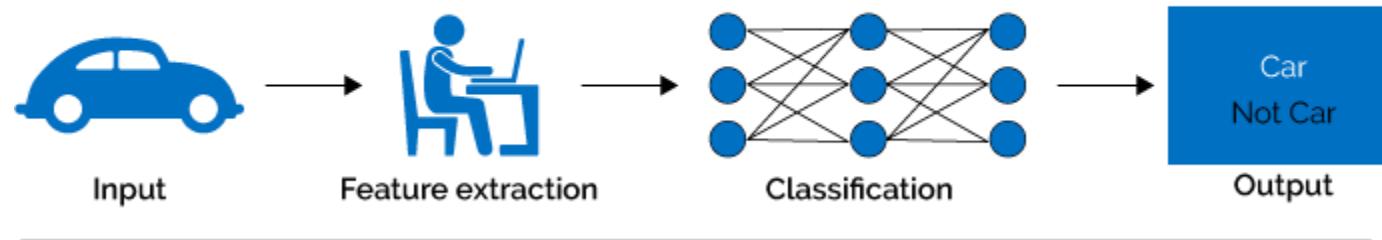
# Graph Representation Learning

Graph Representation Learning alleviates the need to do feature engineering **every single time.**

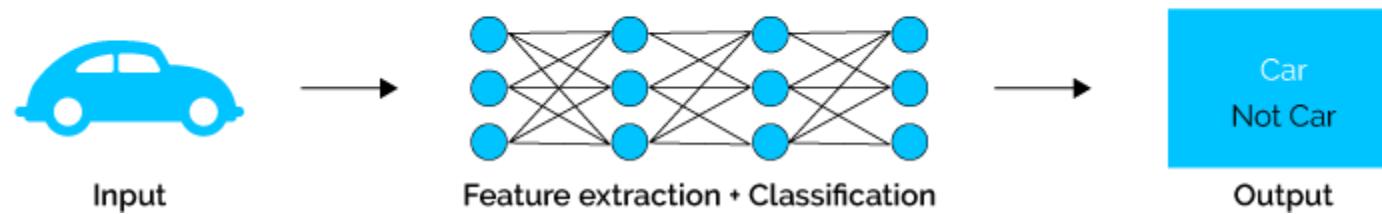


# Recap: Classical ML vs. DL (Better Representations)

## Machine Learning

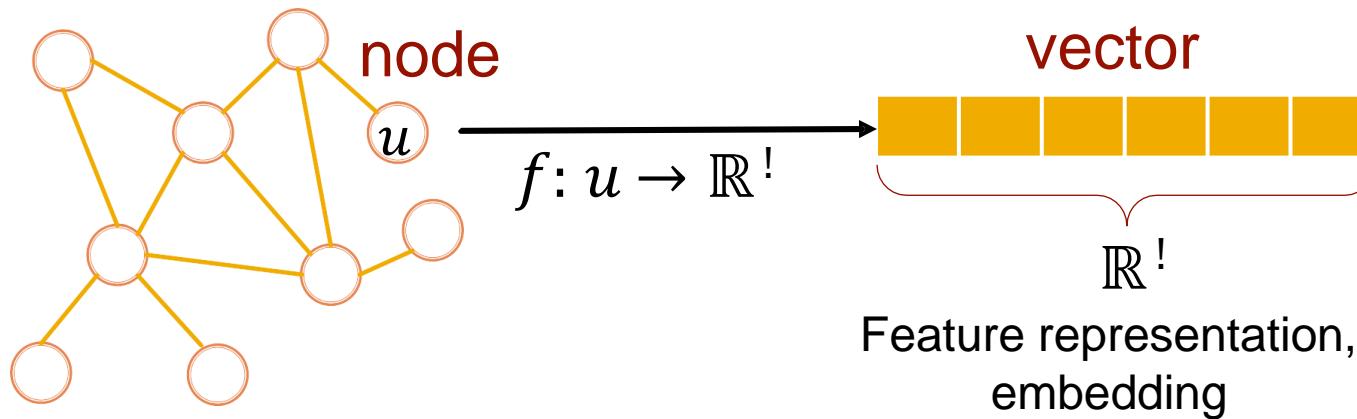


## Deep Learning



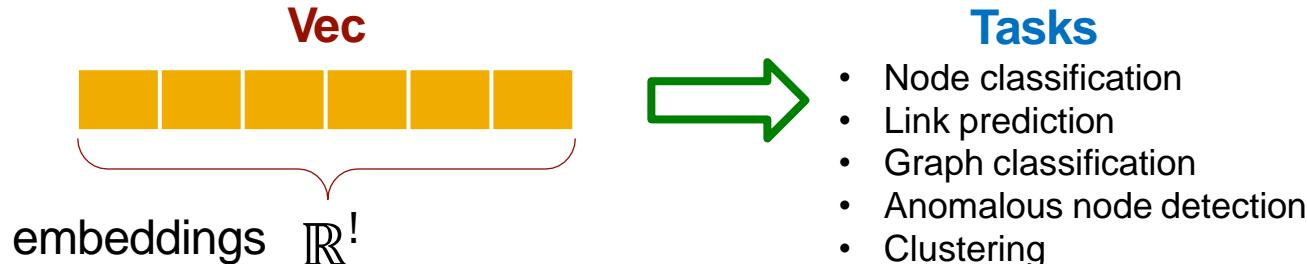
# Graph Representation Learning

**Goal:** Efficient task-independent feature learning for machine learning with graphs!



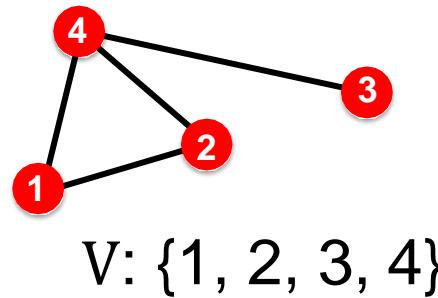
# Why Embedding?

- **Task: Map nodes into an embedding space**
  - Similarity of embeddings between nodes indicates their similarity in the network. For example:
    - Both nodes are close to each other (connected by an edge)
  - Encode network information
  - Potentially used for many downstream predictions



# Setup

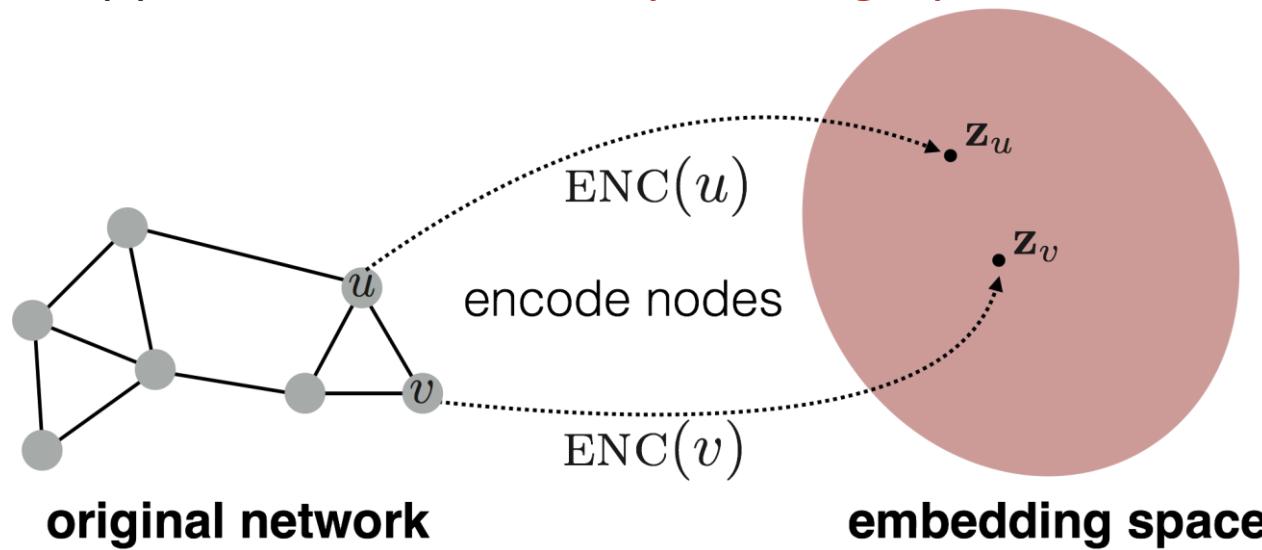
- Assume we have an (undirected) graph  $G$ :
  - $V$  is the vertex set.
  - $A$  is the adjacency matrix (assume binary).
  - **For simplicity: No node features or extra information is used**



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

# Embedding Nodes

- Goal is to encode nodes so that **similarity in the embedding space (e.g., dot product)** approximates **similarity in the graph**



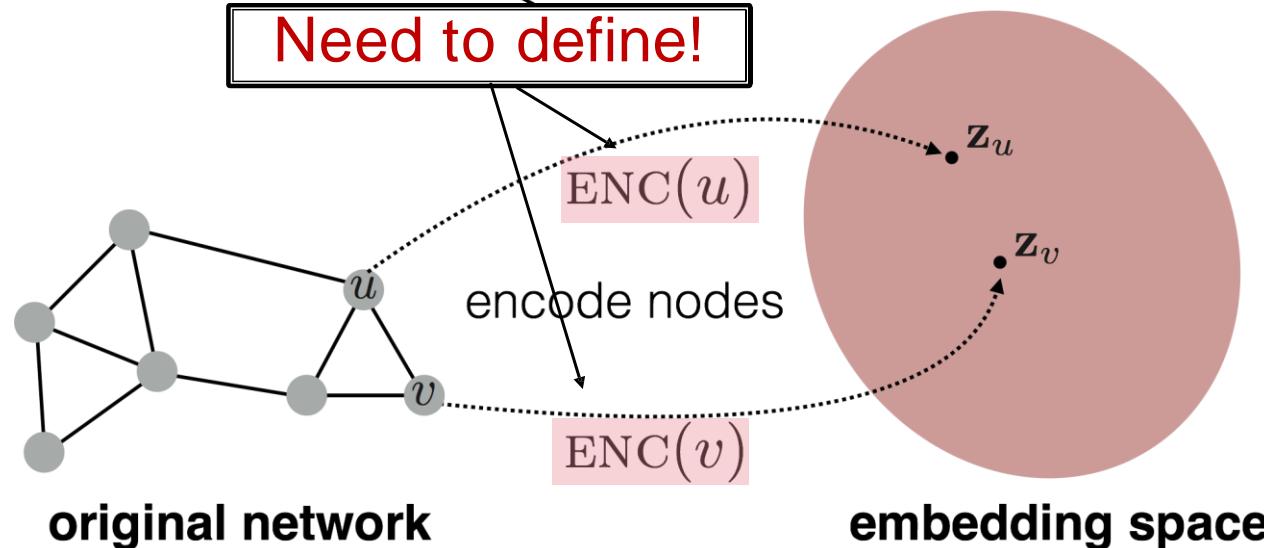
# Embedding Nodes

Goal:  $\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u$

in the original network

Similarity of the embedding

Need to define!



# Learning Node Embeddings

1. **Encoder** maps from nodes to embeddings
2. Define a node similarity function (i.e., a measure of similarity in the original network)
3. **Decoder DEC** maps from embeddings to the similarity score
4. Optimize the parameters of the encoder so that:

DEC( $\mathbf{z}_v^T \mathbf{z}_u$ )

$$\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u$$

in the original network

Similarity of the embedding

# Two Key Components

- **Encoder:** maps each node to a low-dimensional vector

$$\text{ENC}(v) = \mathbf{z}_v \xrightarrow{\text{node in the input graph}} \begin{matrix} d\text{-dimensional} \\ \text{embedding} \end{matrix}$$

- **Similarity function:** specifies how the relationships in vector space map to the relationships in the original network

$$\text{similarity}(u, v) \approx \mathbf{z}_v^T \mathbf{z}_u \xrightarrow{\text{Similarity of } u \text{ and } v \text{ in the original network}} \begin{matrix} \text{Decoder} \\ \text{dot product between node embeddings} \end{matrix}$$

# “Shallow” Encoding

Simplest encoding approach: Encoder is just an embedding-lookup

$$\text{ENC}(v) = \mathbf{z}_v = \mathbf{Z} \cdot v$$

$$\mathbf{Z} \in \mathbb{R}^{d \times |\mathcal{V}|}$$

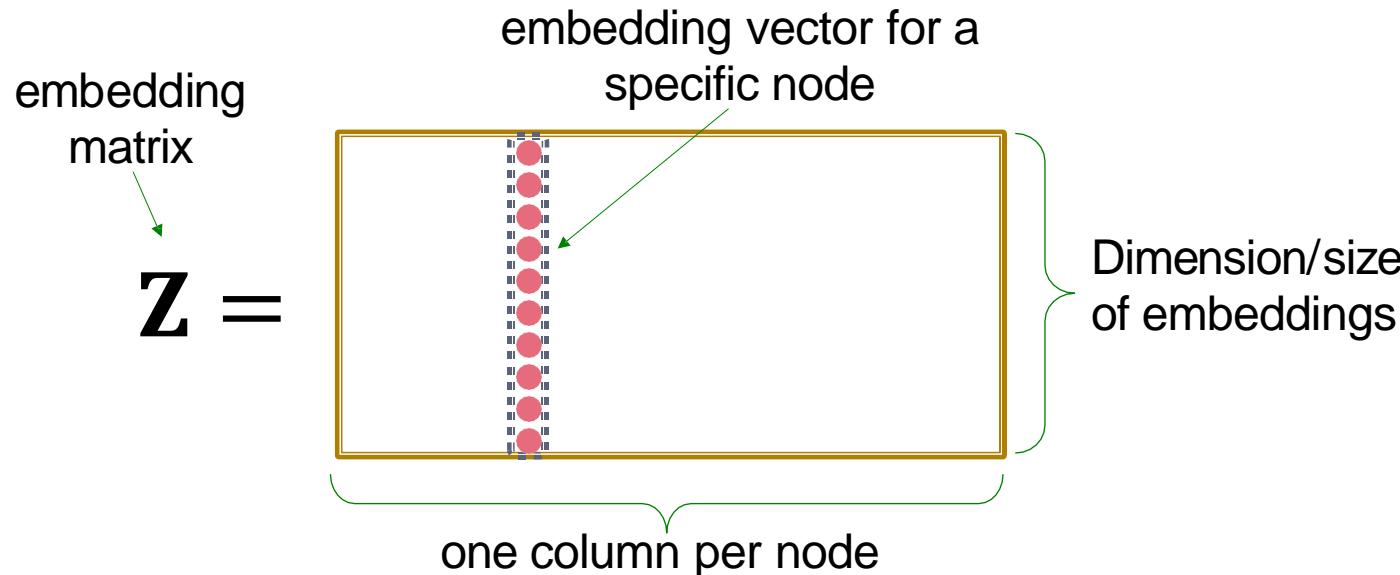
matrix, each column is a node embedding [what we learn / optimize]

$$v \in \mathbb{I}^{|\mathcal{V}|}$$

indicator vector, all zeroes except a one in column indicating node  $v$

# “Shallow” Encoding

Simplest encoding approach: **encoder is just an embedding-lookup**



# “Shallow” Encoding

Simplest encoding approach: **Encoder is just an embedding-lookup**

**Each node is assigned a unique  
embedding vector**

(i.e., we directly optimize  
the embedding of each node)

Many methods: DeepWalk, node2vec

# Framework Summary

- **Encoder + Decoder Framework**
  - Shallow encoder: Embedding lookup
  - Parameters to optimize:  $\mathbf{Z}$  which contains node embeddings  $\mathbf{z}_u$  for all nodes  $u \in V$
  - We will cover deep encoders in the GNNs in 20 mins
- **Decoder:** based on node similarity.
- **Objective:** maximize  $\mathbf{z}_v^T \mathbf{z}_u$  for node pairs  $(u, v)$  that are **similar**

# How to Define Node Similarity?

- Key choice of methods is **how they define node similarity.**
- Should two nodes have a similar embedding if they...
  - are linked?
  - share neighbors?
  - have similar “structural roles”?

# Note on Node Embeddings

- This is **unsupervised/self-supervised** way of learning node embeddings.
  - We are **not** utilizing node labels
  - We are **not** utilizing node features
  - The goal is to directly estimate a set of coordinates (i.e., the embedding) of a node so that some aspect of the network structure is preserved.
- These embeddings are **task independent:**
  - They are not trained for a specific task but can be used for any task.

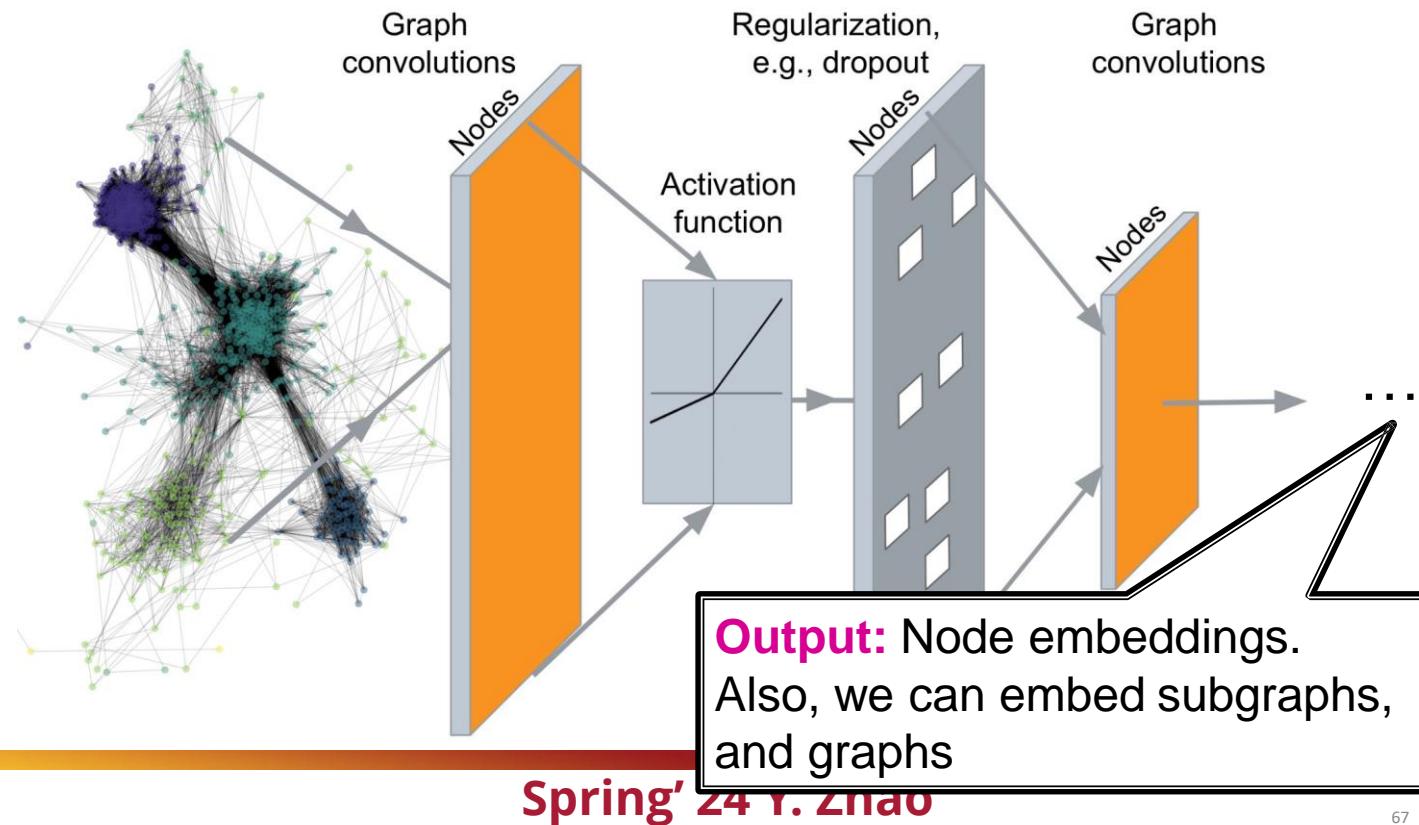
# Today: Deep Graph Encoders

- **Now:** discuss deep learning methods based on **graph neural networks (GNNs)**:

$\text{ENC}(\nu) =$  **multiple layers of  
non-linear transformations  
based on graph structure**

- **Note:** All these deep encoders can be **combined with node similarity functions** defined earlier.

# Deep Graph Encoders

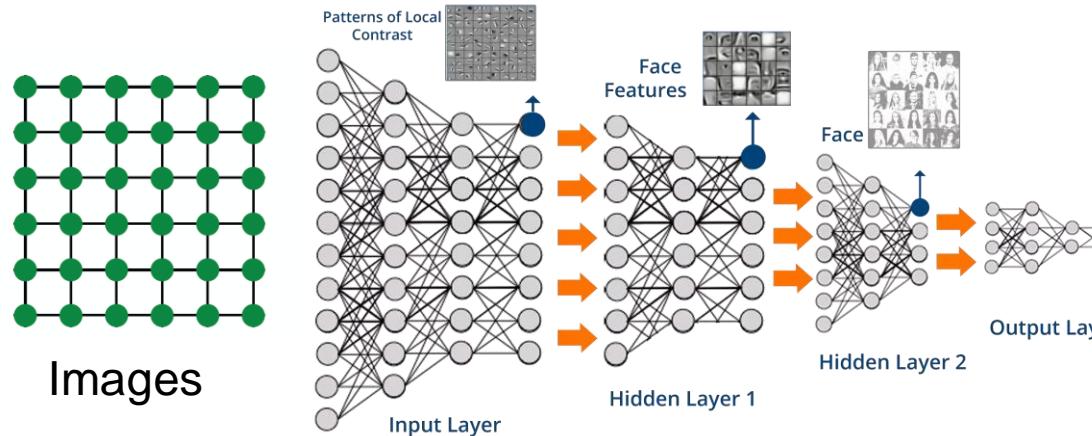


# Tasks on Networks

## Tasks we will be able to solve:

- Node classification
  - Predict the type of a given node
- Link prediction
  - Predict whether two nodes are linked
- Community detection
  - Identify densely linked clusters of nodes
- Network similarity
  - How similar are two (sub)networks

# Modern ML Toolbox

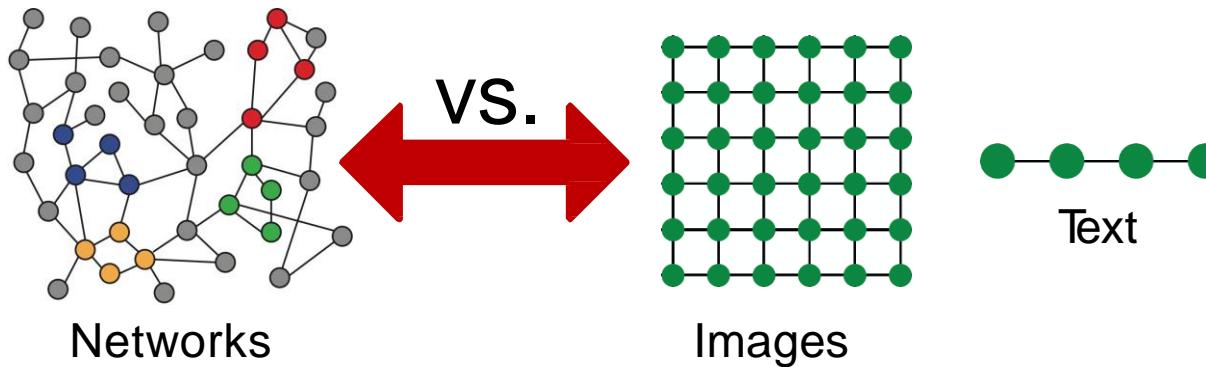


Modern deep learning toolbox is designed  
for simple sequences & grids

# Why is it Hard?

**But networks are far more complex!**

- Arbitrary size and complex topological structure (i.e., no spatial locality like grids)



- No fixed node ordering or reference point
- Often dynamic and have multimodal features

# Summary: Basics of Deep Learning

- **Loss function:**

$$\min_{\Theta} \mathcal{L}(y, f_{\Theta}(x))$$

- $f$  can be a simple linear layer, an MLP, or other neural networks (e.g., a GNN later)
- Sample a minibatch of input  $x$
- **Forward propagation:** Compute  $\mathcal{L}$  given  $x$
- **Back-propagation:** Obtain gradient  $\nabla_{\Theta} \mathcal{L}$  using a chain rule.
- Use **stochastic gradient descent (SGD)** to optimize  $\mathcal{L}$  for  $\Theta$  over many iterations.

# Content

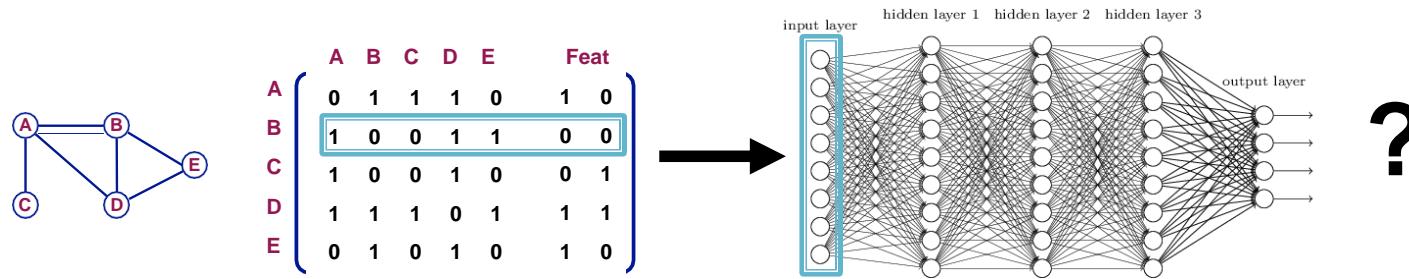
- **Local network neighborhoods:**
  - Describe aggregation strategies
  - Define computation graphs
- **Stacking multiple layers:**
  - Describe the model, parameters, training
  - How to fit the model?
  - Simple example for unsupervised and supervised training

# Setup

- Assume we have a graph  $G$ :
  - $V$  is the **vertex set**
  - $A$  is the **adjacency matrix** (assume binary)
  - $X \in \mathbb{R}^{|V| \times m}$  is a matrix of **node features**
  - $v$ : a node in  $V$ ;  $N(v)$ : 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:  $[1, 1, \dots, 1]$

# A Naïve Approach

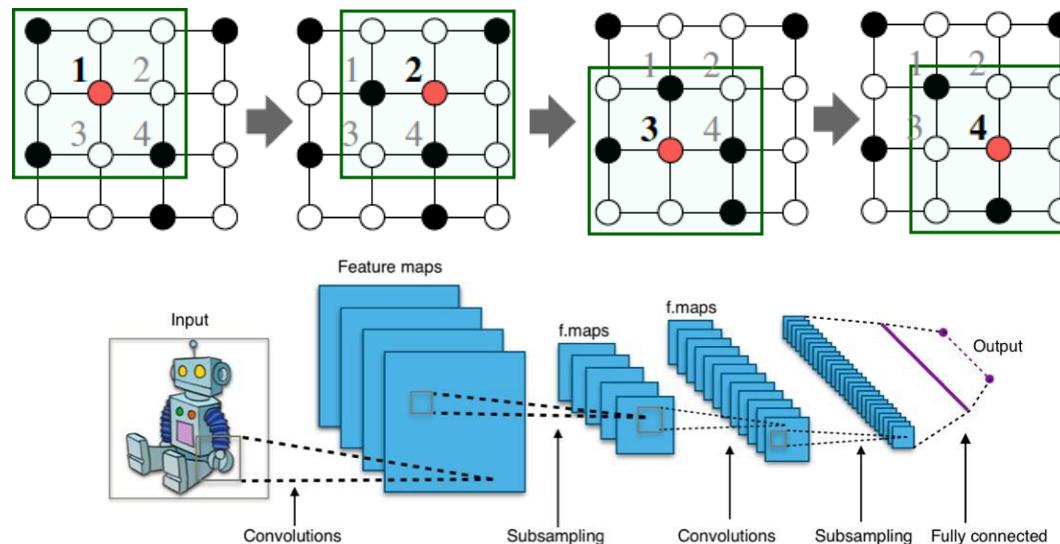
- Join adjacency matrix and features
- Feed them into a deep neural net:



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

# Idea: Convolutional Networks

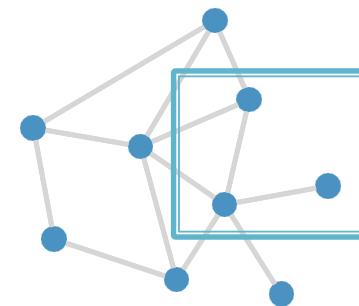
CNN on an image:



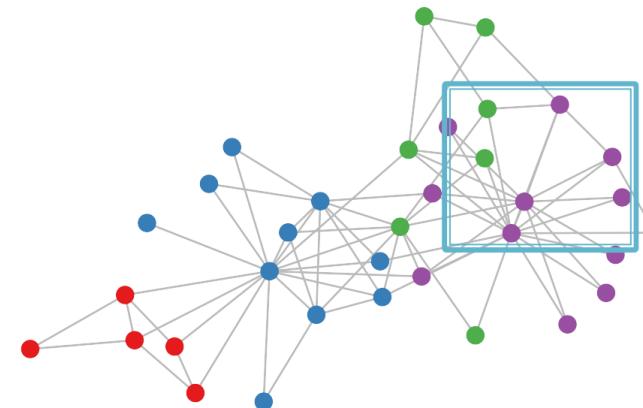
Goal is to generalize convolutions beyond simple lattices  
 Leverage node features/attributes (e.g., text, images)

# Real-World Graphs

But our graphs look like this:



or this:



- There is no fixed notion of locality or sliding window on the graph
- Graph is permutation invariant

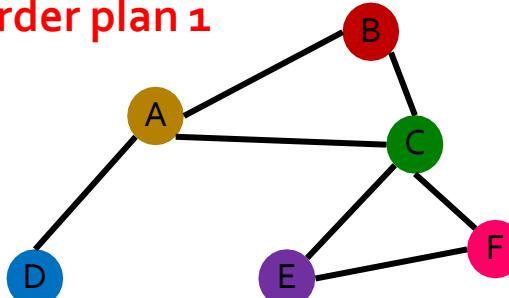
# Permutation Invariance

- **Graph does not have a canonical order of the nodes!**
- We can have many different order plans.

# Permutation Invariance

- Graph does not have a canonical order of the nodes!

Order plan 1



Node features  $X_1$

A	Dark Brown
B	Red
C	Green
D	Blue
E	Purple
F	Pink

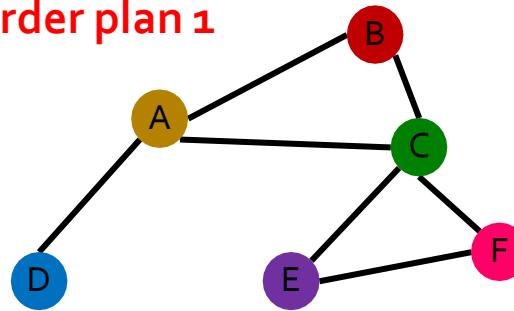
Adjacency matrix  $A_1$

A	B	C	D	E	F
A	Gray	Light Blue	Light Blue	Gray	Gray
B	Light Blue	Gray	Light Blue	Light Blue	Gray
C	Light Blue	Gray	Gray	Light Blue	Light Blue
D	Light Blue	Gray	Gray	Gray	Gray
E	Gray	Gray	Light Blue	Gray	Light Blue
F	Gray	Light Blue	Light Blue	Gray	Gray

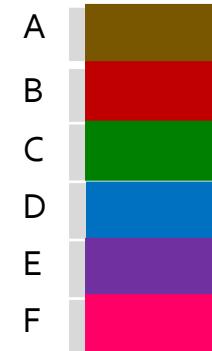
# Permutation Invariance

- Graph does not have a canonical order of the nodes!

Order plan 1



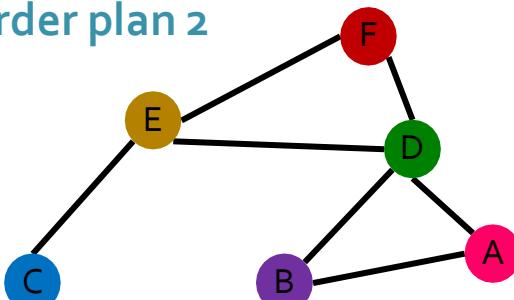
Node features  $X_1$



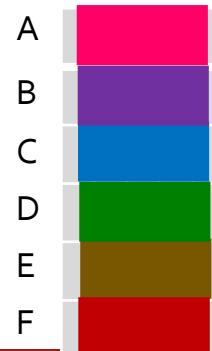
Adjacency matrix  $A_1$

	A	B	C	D	E	F
A	Gray	Blue	Blue	Blue	Gray	Gray
B	Blue	Gray	Blue	Gray	Gray	Gray
C	Blue	Blue	Gray	Blue	Gray	Blue
D	Blue	Gray	Gray	Gray	Gray	Blue
E	Gray	Gray	Blue	Gray	Gray	Blue
F	Gray	Gray	Gray	Blue	Blue	Gray

Order plan 2



Node features  $X_2$



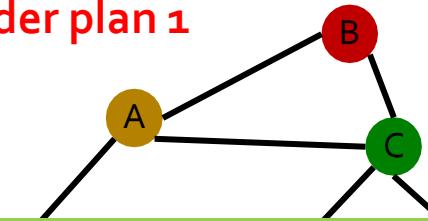
Adjacency matrix  $A_2$

	A	B	C	D	E	F
A	Gray	Blue	Blue	Blue	Gray	Gray
B	Blue	Gray	Blue	Gray	Gray	Gray
C	Blue	Blue	Gray	Blue	Gray	Blue
D	Blue	Gray	Gray	Gray	Gray	Blue
E	Gray	Gray	Blue	Gray	Gray	Blue
F	Gray	Gray	Gray	Blue	Blue	Gray

# Permutation Invariance

- Graph does not have a canonical order of the nodes!

Order plan 1



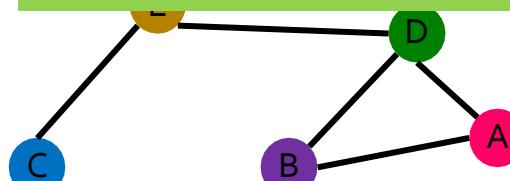
Node features  $X_1$

A					
B					
C					
D					

Adjacency matrix  $A_1$

A	B	C	D	E	F
A	1	1	0	0	0
B	0	1	1	0	0
C	0	0	1	1	0
D	0	0	0	1	0

Order plan 2



C					
D					
E					
F					

B	C	D	E	F
B	1	0	0	0
C	0	1	0	0
D	0	0	1	0
E	0	0	0	1
F	0	0	0	0

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

# Permutation Invariance

What does it mean by “graph representation is same for two order plans”?

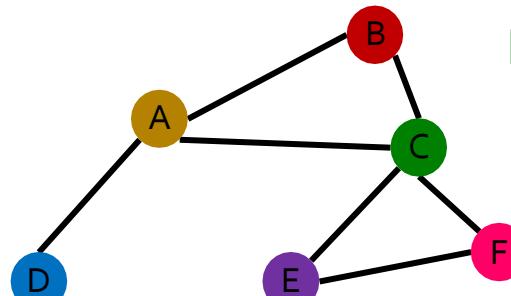
- Consider we learn a function  $f$  that maps a graph  $G = (A, X)$  to a vector  $\mathbb{R}^d$  then

$$f(A_1, X_1) = f(A_2, X_2)$$

In other words,  $f$  maps a graph to a  $d$ -dim embedding

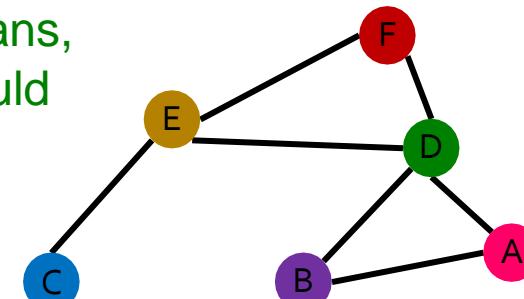
$A$  is the adjacency matrix  
 $X$  is the node feature matrix

Order plan 1:  $A_1, X_1$



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

Order plan 2:  $A_2, X_2$



# Permutation Invariance

**What does it mean by “graph representation is same for two order plans”?**

- Consider we learn a function  $f$  that maps a graph  $G = (A, X)$  to a vector  $\mathbb{R}^d$ .

$A$  is the adjacency matrix  
 $X$  is the node feature matrix

- Then, if  $f(A_i, X_i) = f(A_j, X_j)$  for any order plan  $i$  and  $j$ , we formally say  $f$  is a **permutation invariant**

**function.** For a graph with  $|V|$  nodes, there are  $|V|!$  different order plans.

$m$ ... each node has a  $m$ -dim feature vector associated with it.

- Definition:** For any graph function  $f: \mathbb{R}^{|V| \times m} \times \mathbb{R}^{|V| \times |V|} \rightarrow \mathbb{R}^d$ ,  $f$  is **permutation-invariant** if  $f(A, X) = f(PAP^T, PX)$  for any permutation  $P$ .

$d$ ... output embedding dimensionality of embedding the graph  $G = (A, X)$

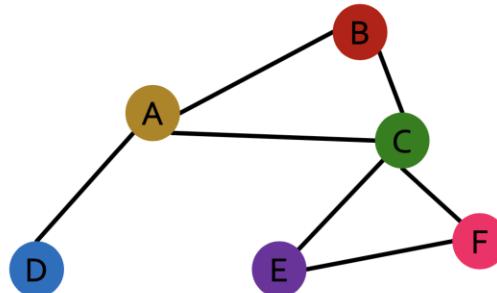
Permutation  $P$ : a shuffle of the node order  
Example: (A,B,C)->(B,C,A)

# Permutation Equivariance

**For node representation:** We learn a function  $f$  that maps nodes of  $G$  to a matrix  $\mathbb{R}^{|V| \times d}$ .

In other words, each node in  $V$  is mapped to a  $d$ -dim embedding.

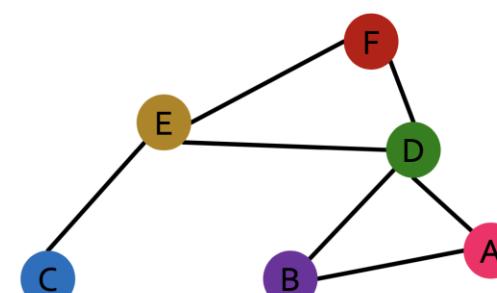
Order plan 1:  $A_1, X_1$



$$f(A_1, X_1) = \begin{matrix} & \text{A} & \text{B} \\ \text{A} & \text{---} & \text{---} \\ \text{B} & \text{---} & \text{---} \\ \text{C} & \text{---} & \text{---} \\ \text{D} & \text{---} & \text{---} \\ \text{E} & \text{---} & \text{---} \\ \text{F} & \text{---} & \text{---} \end{matrix}$$

The matrix shows that node A has two yellow entries, B has two red entries, C has two green entries, D has two blue entries, E has two purple entries, and F has two orange-red entries.

Order plan 2:  $A_2, X_2$



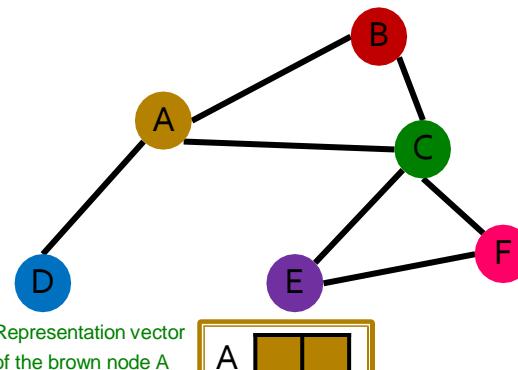
$$f(A_2, X_2) = \begin{matrix} & \text{A} & \text{B} \\ \text{A} & \text{---} & \text{---} \\ \text{B} & \text{---} & \text{---} \\ \text{C} & \text{---} & \text{---} \\ \text{D} & \text{---} & \text{---} \\ \text{E} & \text{---} & \text{---} \\ \text{F} & \text{---} & \text{---} \end{matrix}$$

The matrix shows that node A has two red entries, B has two purple entries, C has two blue entries, D has two green entries, E has two yellow entries, and F has two orange-red entries.

# Permutation Equivariance

**For node representation:** We learn a function  $f$  that maps nodes of  $G$  to a matrix  $\mathbb{R}^{|v| \times d}$ .

Order plan 1:  $A_1, X_1$

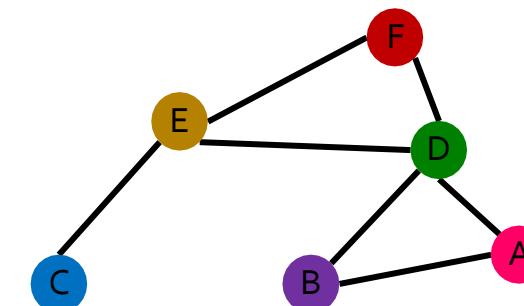


$$f(A_1, X_1) =$$

A	B	C	D	E	F

For two order plans, the vector of node at the same position in the graph is the same!

Order plan 2:  $A_2, X_2$



$$f(A_2, X_2) =$$

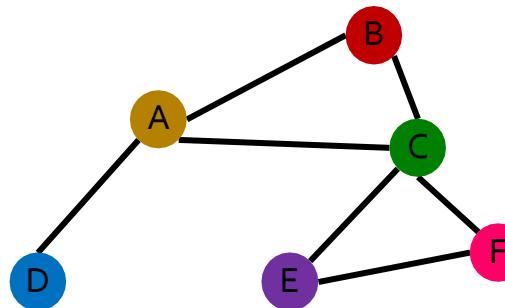
A	B	C	D	E	F

Representation vector of the brown node E

# Permutation Equivariance

**For node representation:** We learn a function  $f$  that maps nodes of  $G$  to a matrix  $\mathbb{R}^{|v| \times d}$ .

Order plan 1:  $A_1, X_1$

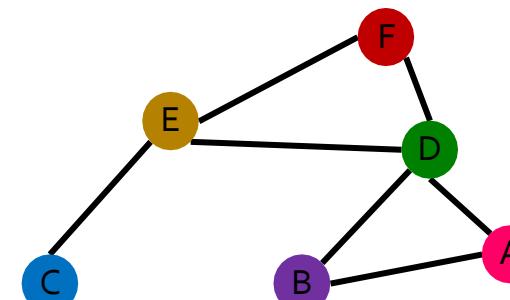


$$f(A_1, X_1) = \begin{matrix} & \text{A} & \text{B} & \text{C} & \text{D} & \text{E} \\ \text{A} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{B} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{C} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{D} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{E} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{F} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \end{matrix}$$

Representation vector of the green node C

For two order plans, the vector of node at the same position in the graph is the same!

Order plan 2:  $A_2, X_2$



$$f(A_2, X_2) = \begin{matrix} & \text{A} & \text{B} & \text{C} & \text{D} & \text{E} & \text{F} \\ \text{A} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{B} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{C} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{D} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{E} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \\ \text{F} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} & \text{---} \end{matrix}$$

Representation vector of the green node D

# Permutation Equivariance

**For node representation:**

- Consider we learn a function  $f$  that maps a graph  $G = (A, X)$  to a matrix  $\mathbb{R}^{|V| \times d}$
- If the output vector of a node at the same position in the graph remains unchanged for any order plan, we say  $f$  is **permutation equivariant**.
- **Definition:** For any **node** function  $f: \mathbb{R}^{|V| \times m} \times \mathbb{R}^{|V| \times |V|} \rightarrow \mathbb{R}^{|V| \times d}$ ,  $f$  is **permutation-equivariant** if  $Pf(A, X) = f(PAP^T, PX)$  for any permutation  $P$ .  $f$  maps each node in  $V$  to a  $d$ -dim embedding.

$m$ ... each node has a  $m$ -dim feature vector associated with it.

# Summary: Invariance and Equivariance

## ■ Permutation-invariant

$$f(A, X) = f(PAP^T, PX)$$

Permute the input, the output  
stays the same.

(map a graph to a vector)

graph level

## ■ Permutation-equivariant

$$Pf(A, X) = f(PAP^T, PX)$$

Permute the input, output also  
permutes accordingly.  
(map a graph to a matrix)

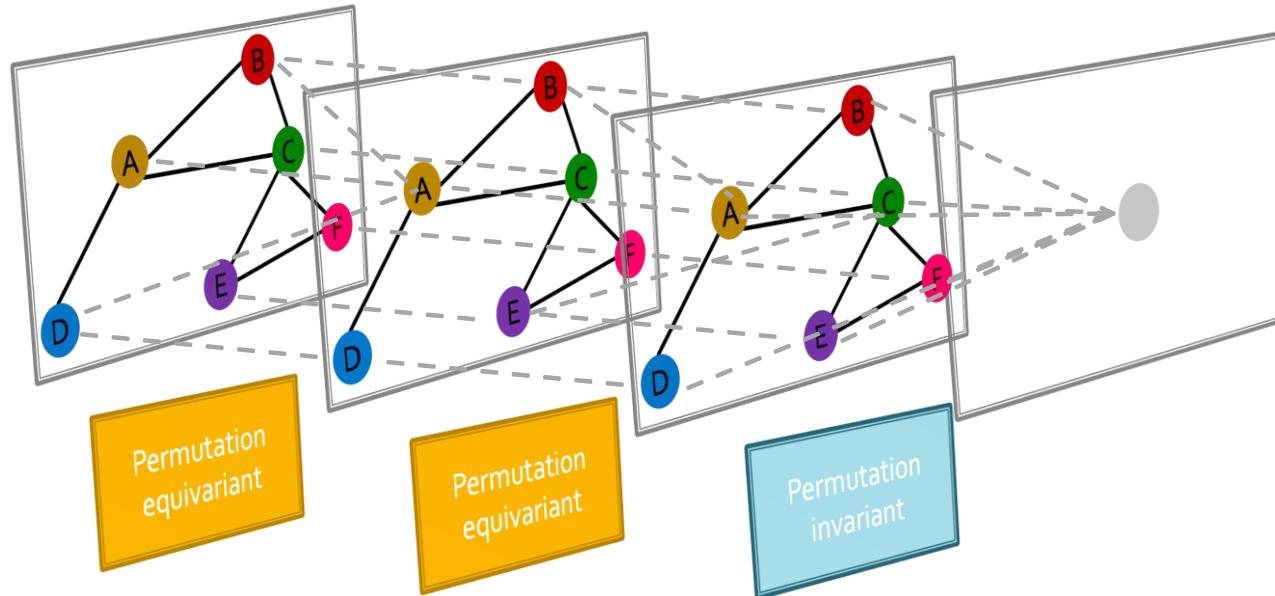
node level

## ■ Examples:

- $f(A, X) = \mathbf{1}^T \mathbf{X}$  : Permutation-**invariant**
  - Reason:  $f(PAP^T, PX) = \mathbf{1}^T \mathbf{P} \mathbf{X} = \mathbf{1}^T \mathbf{X} = f(A, X)$
- $f(A, X) = \mathbf{X}$  : Permutation-**equivariant**
  - Reason:  $f(PAP^T, PX) = \mathbf{P} \mathbf{X} = Pf(A, X)$
- $f(A, X) = \mathbf{A} \mathbf{X}$  : Permutation-**equivariant**
  - Reason:  $f(PAP^T, PX) = PAP^T \mathbf{P} \mathbf{X} = \mathbf{P} \mathbf{A} \mathbf{X} = Pf(A, X)$

# Graph Neural Network Overview

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

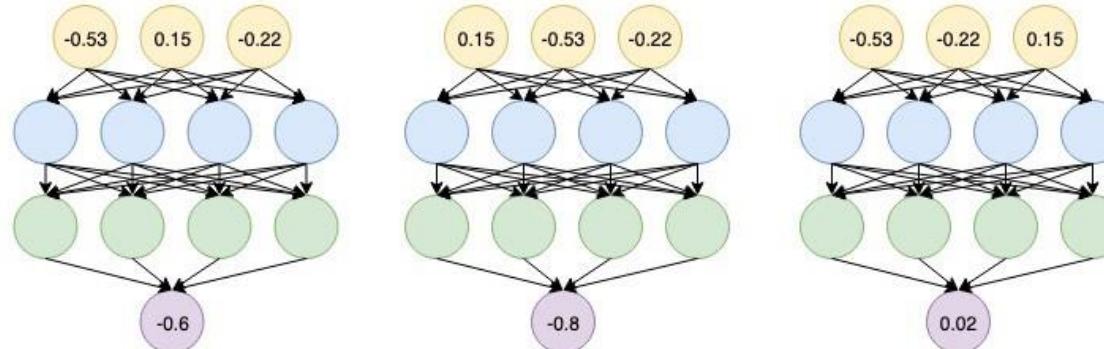


# Graph Neural Network Overview

Are other neural network architectures, e.g.,  
MLPs, permutation invariant / equivariant?

- No.

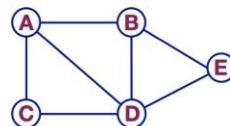
Switching the order of the  
input leads to different  
outputs!



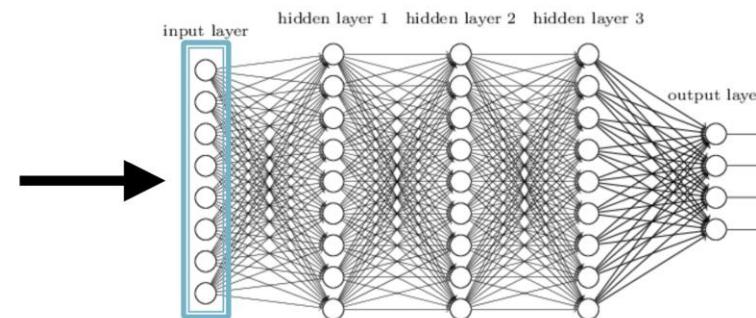
# Graph Neural Network Overview

Are other neural network architectures, e.g.,  
MLPs, permutation invariant / equivariant?

- No.



	A	B	C	D	E	Feat
A	0	1	1	1	0	1 0
B	1	0	0	1	1	0 0
C	1	0	0	1	0	0 1
D	1	1	1	0	1	1 1
E	0	1	0	1	0	1 0



This explains why **the naïve MLP approach fails for graphs!**

# Graph Neural Network Overview

- Are any neural network architectures, e.g.,

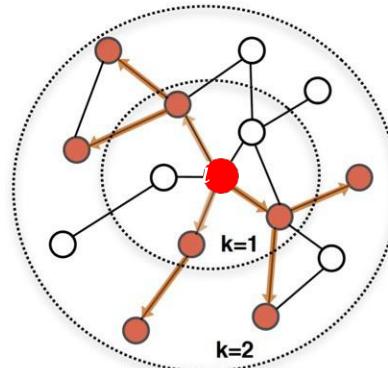


Next: Design graph neural networks that are permutation invariant / equivariant by passing and aggregating information from neighbors!

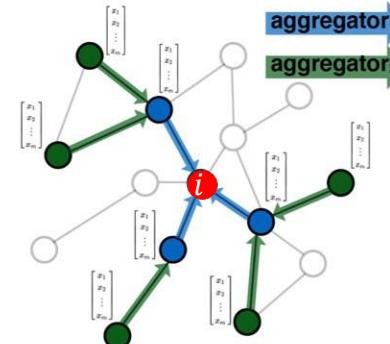
?

# Graph Convolutional Networks

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



Determine node computation graph

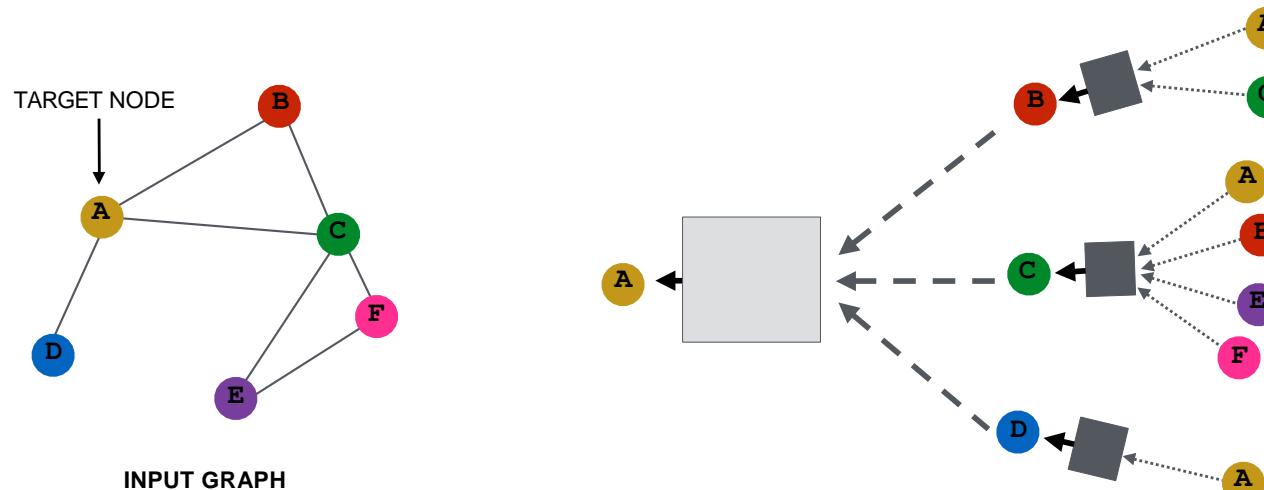


Propagate and transform information

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

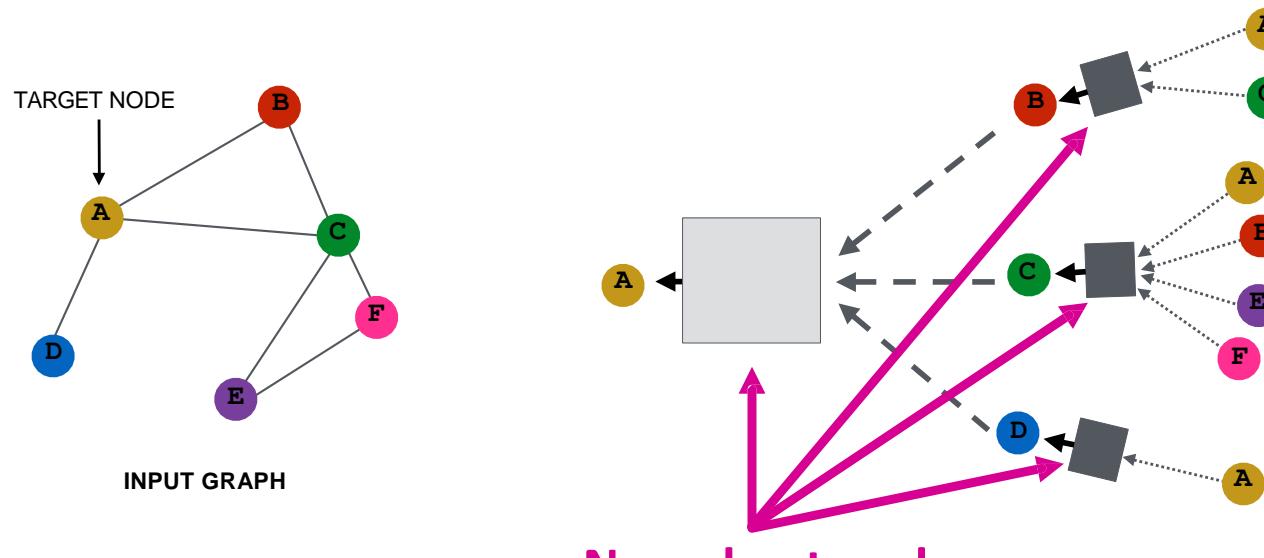
# Idea: Aggregate Neighbors

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



# Idea: Aggregate Neighbors

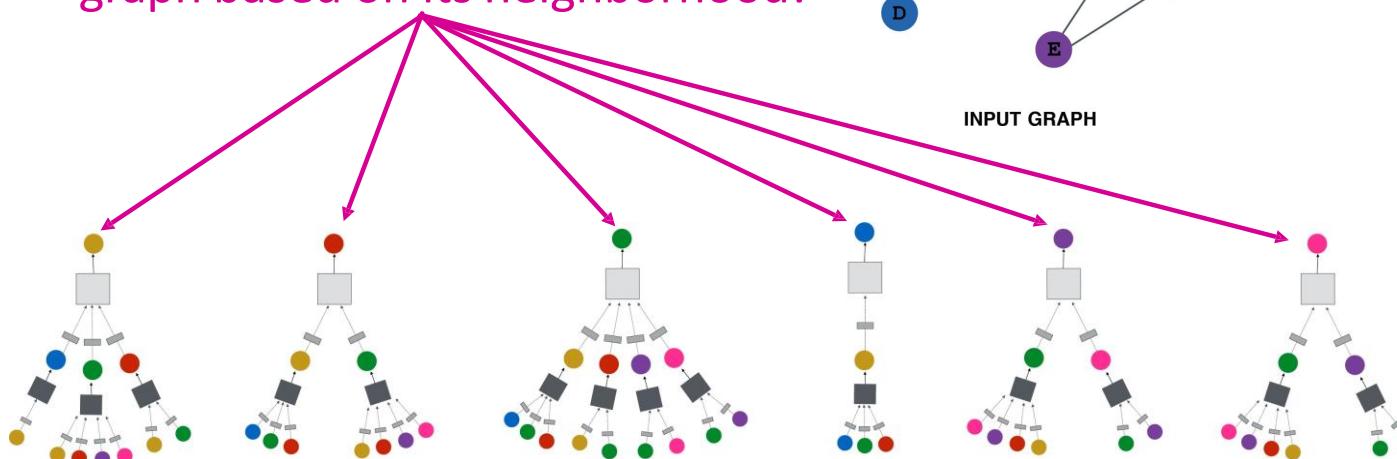
- **Intuition:** Nodes aggregate information from their neighbors using neural networks



# Idea: Aggregate Neighbors

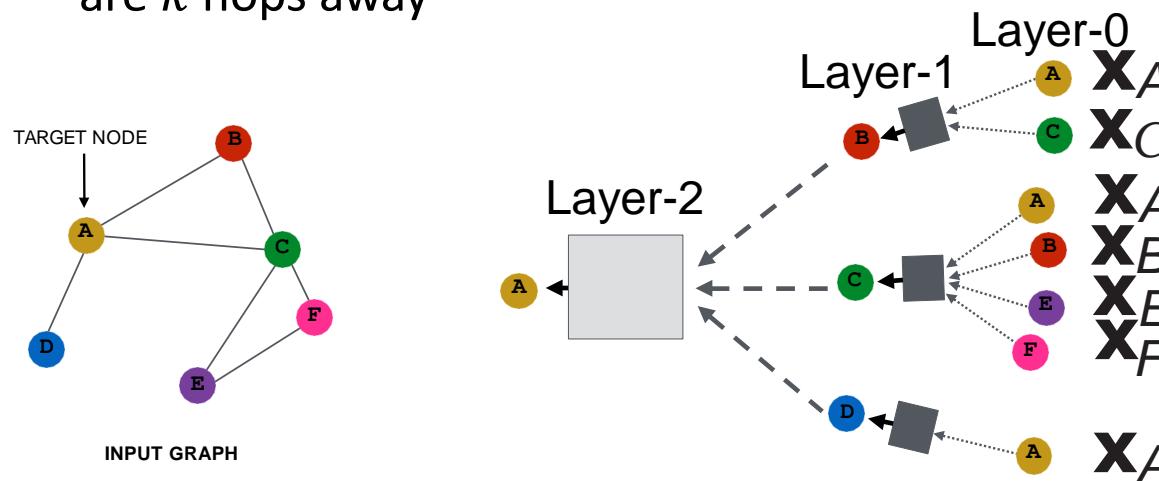
- **Intuition:** Network neighborhood defines a computation graph

Every node defines a computation graph based on its neighborhood!



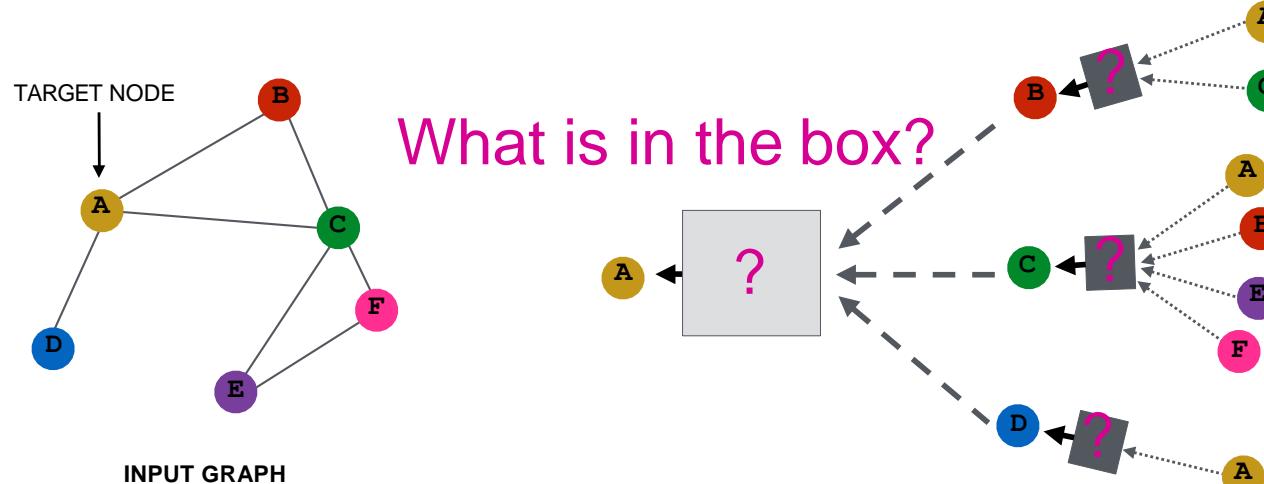
# Deep Model: Many Layers

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



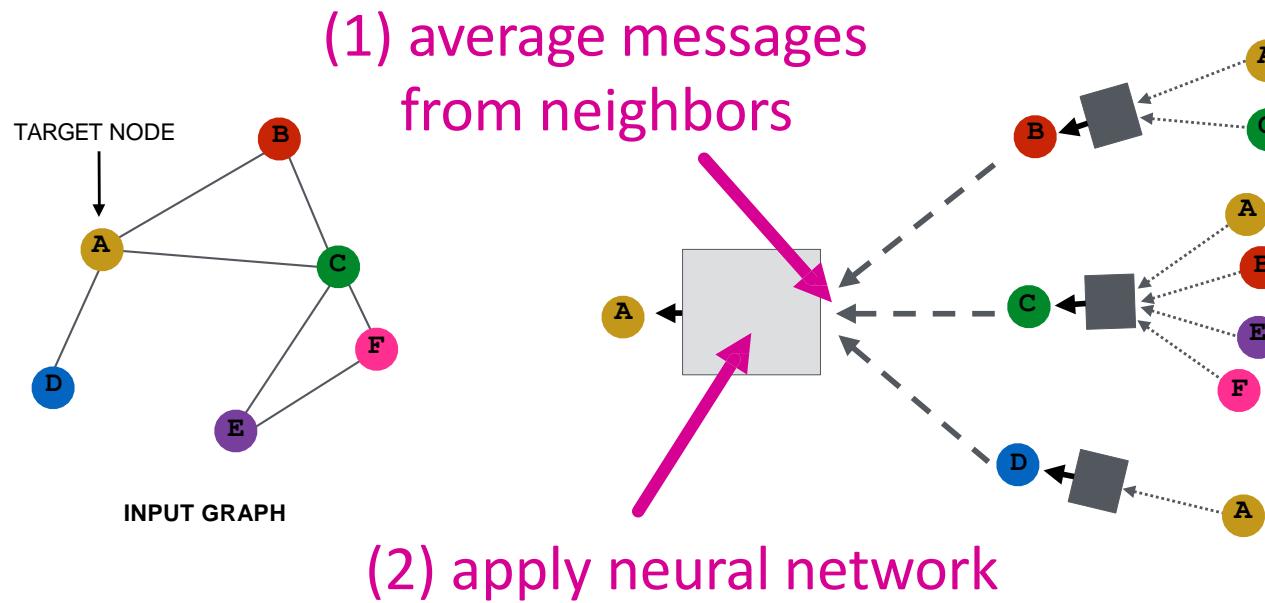
# Neighborhood Aggregation

- **Neighborhood aggregation:** Key distinctions are in how different approaches aggregate information across the layers



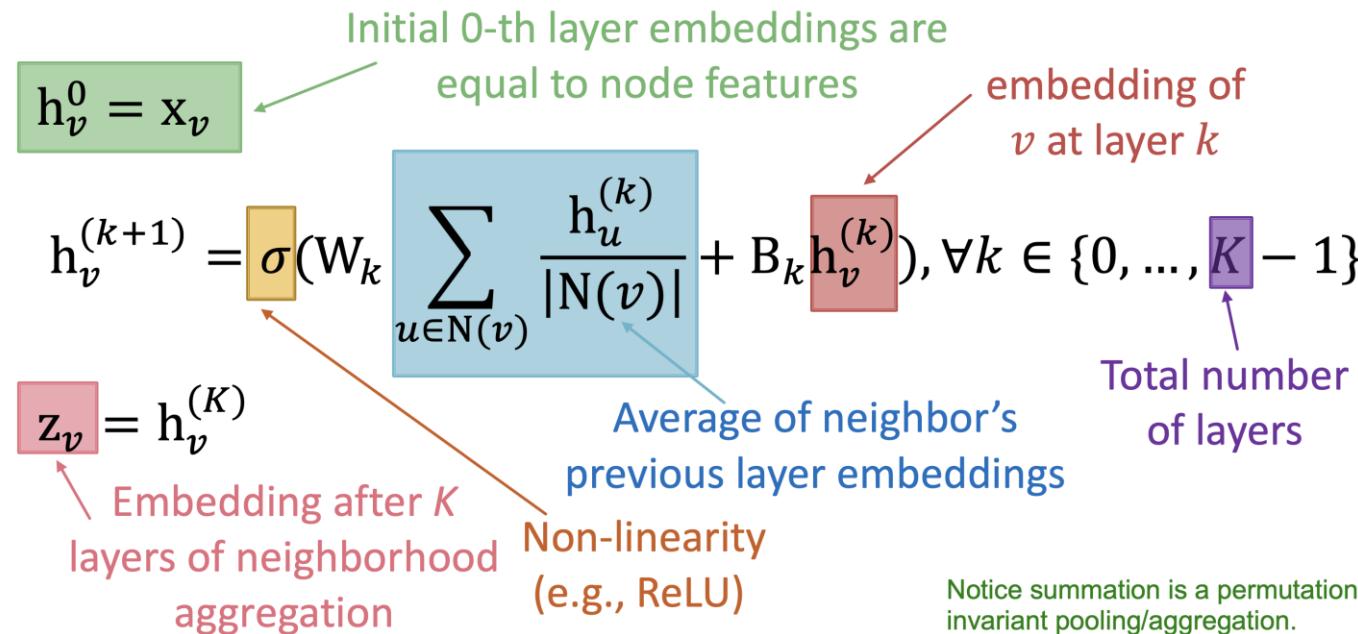
# Neighborhood Aggregation

- **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

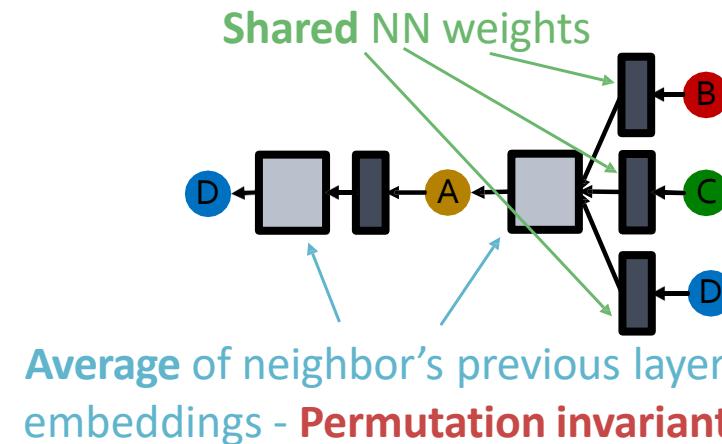
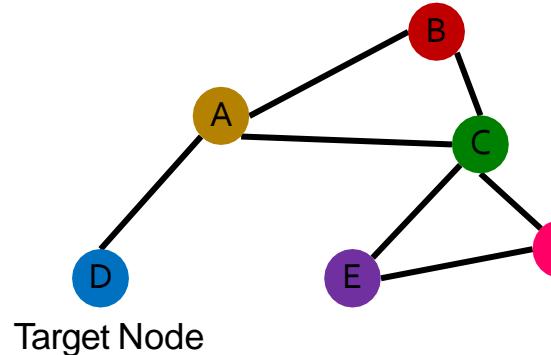


Notice summation is a permutation invariant pooling/aggregation.

# GCN: Invariance and Equivariance

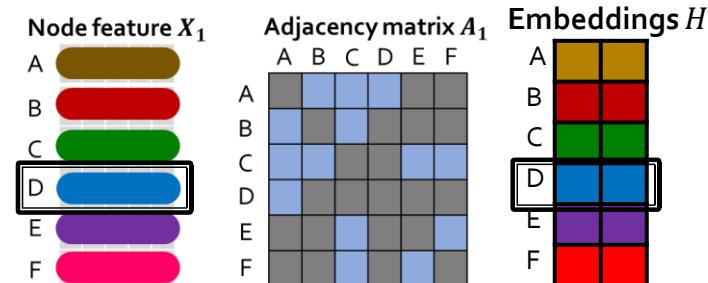
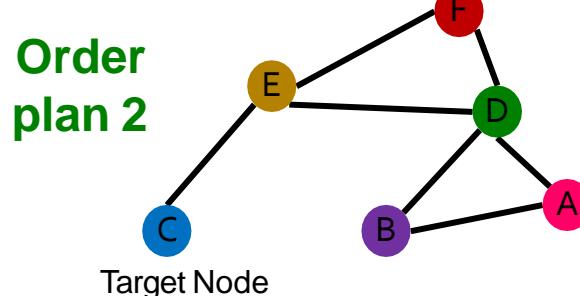
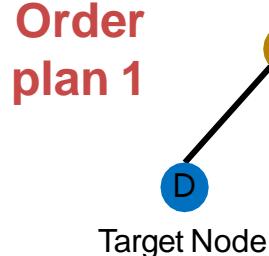
What are the **invariance** and **equivariance** properties for a GCN?

- Given a **node**, the GCN that computes its embedding is **permutation invariant**

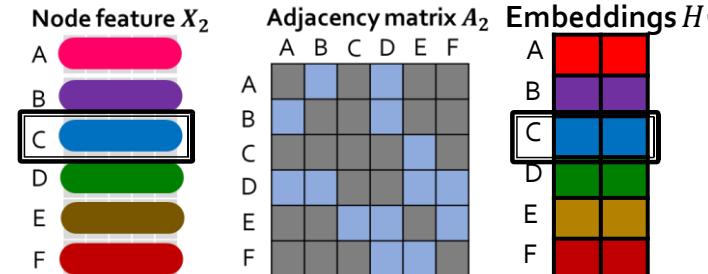


# GCN: Invariance and Equivariance

- Considering all nodes in a graph, GCN computation is **permutation equivariant**



Permute the input, the output also permutes accordingly - **permutation equivariant**



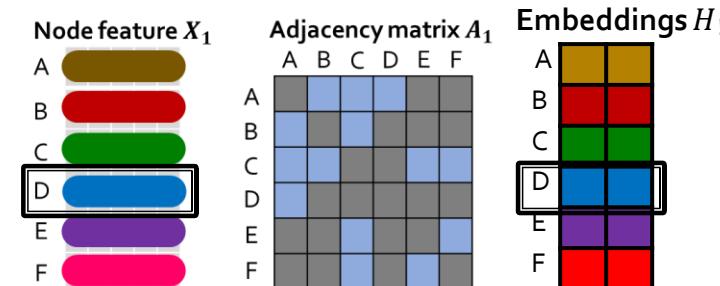
# GCN: Invariance and Equivariance

- Considering all nodes in a graph, GCN computation is **permutation equivariant**

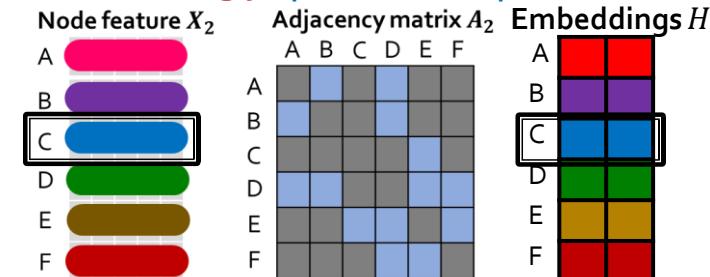
## Detailed reasoning:

- The rows of **input node features** and **output embeddings** are aligned
- We know computing the embedding of a **given node** with GCN is **invariant**.
- So, after permutation, the **location of a given node in the input node feature matrix** is changed, and the **the output embedding of a given node stays the same** (the colors of node feature and embedding are **matched**)

This is **permutation equivariant**

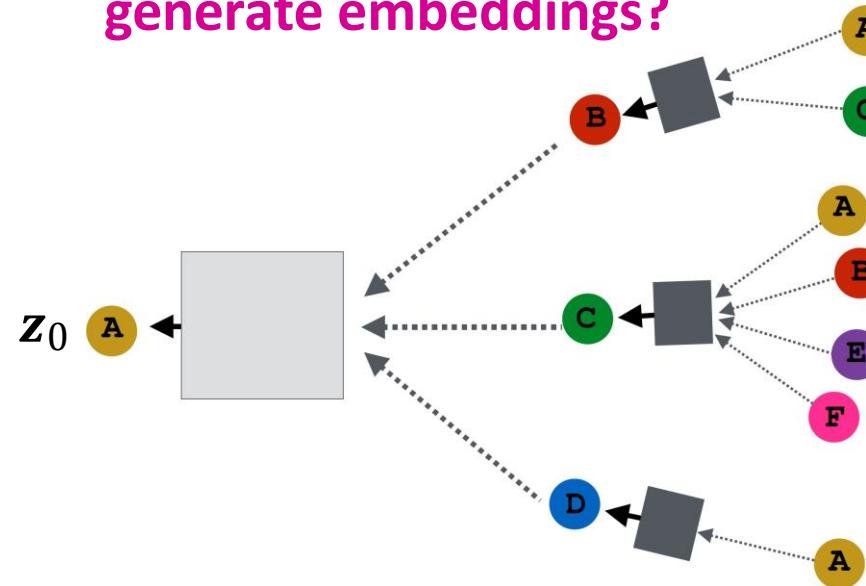


Permute the input, the output also permutes accordingly - **permutation equivariant**



# Training the Model

How do we train the GCN to generate embeddings?



Need to define a loss function on the embeddings.

# Model Parameters

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

$$\begin{aligned}
 h_v^{(0)} &= x_v \\
 h_v^{(k+1)} &= \sigma\left(W_k \sum_{u \in N(v)} \frac{h_u^{(k)}}{|N(v)|} + B_k h_v^{(k)}\right), \forall k \in \{0..K-1\} \\
 z_v &= h_v^{(K)}
 \end{aligned}$$

Final node embedding

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

$h_v^k$ : the hidden representation of node  $v$  at layer  $k$

- $W_k$ : weight matrix for neighborhood aggregation
- $B_k$ : weight matrix for transforming hidden vector of self

# How to Train A GNN

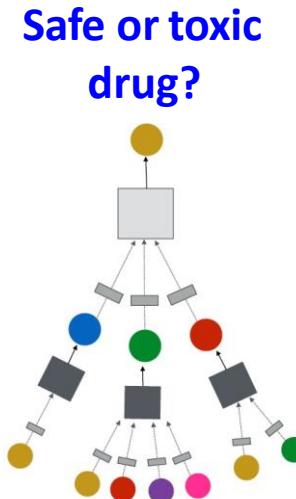
- Node embedding  $\mathbf{z}_0$  is a function of input graph
- **Supervised setting:** We want to minimize loss  $\mathcal{L}$ :

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

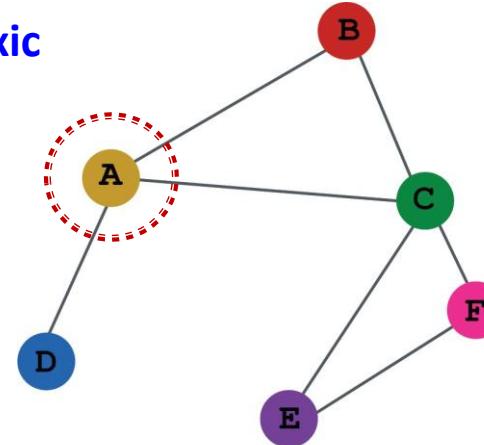
- $\mathbf{y}$ : node label
- $\mathcal{L}$  could be L2 if  $\mathbf{y}$  is real number, or cross entropy if  $\mathbf{y}$  is categorical
- **Unsupervised setting:**
  - No node label available
  - **Use the graph structure as the supervision!**

# Supervised Training

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



Safe or toxic drug?

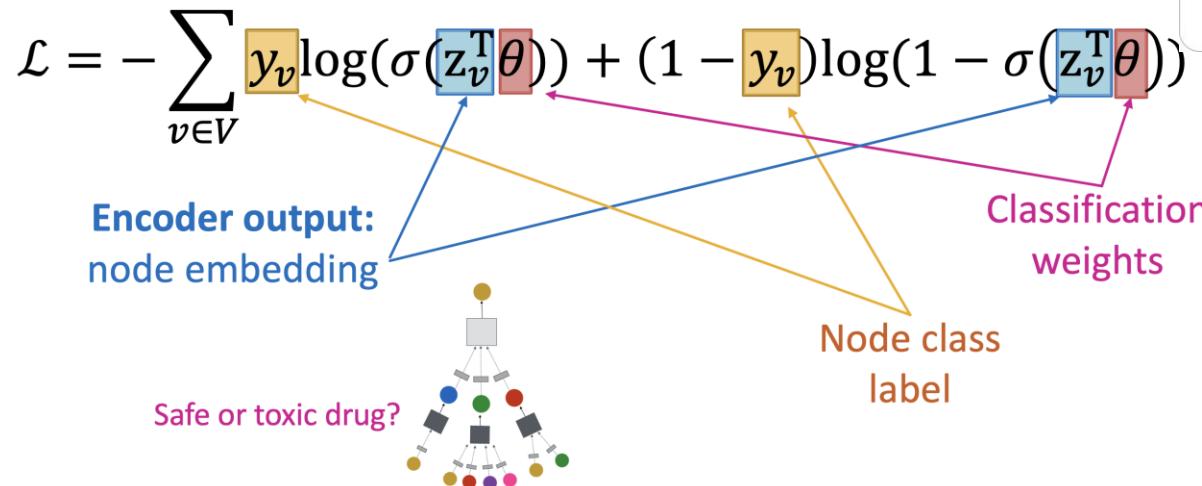


E.g., a drug-drug interaction network

# Supervised Training

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

- Use cross entropy loss



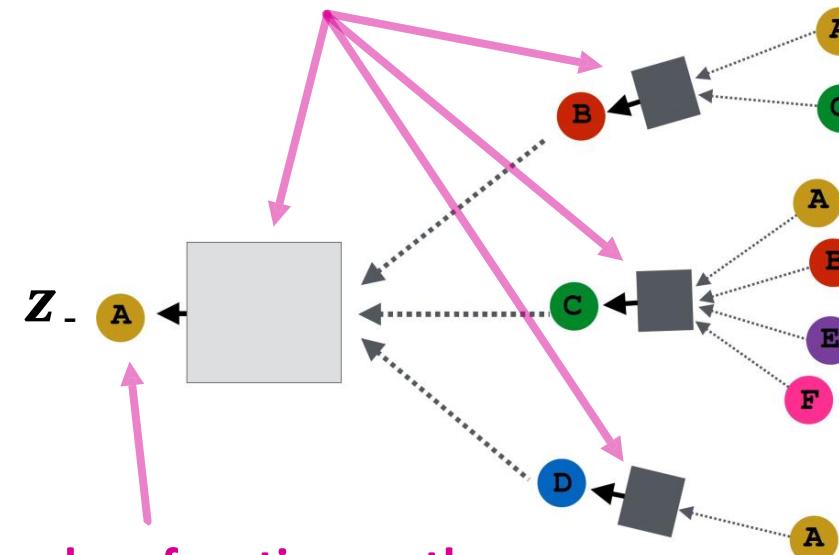
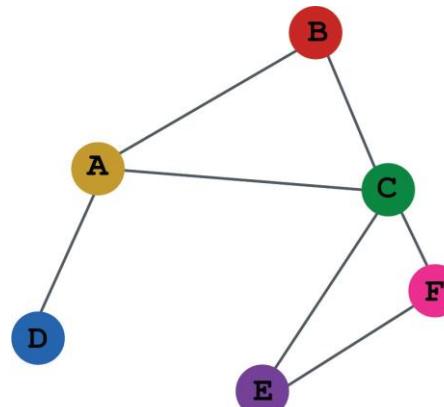
True probability distribution (one-shot)

$$H(p, q) = - \sum_{x \in \text{classes}} p(x) \log q(x)$$

Your model's predicted probability distribution

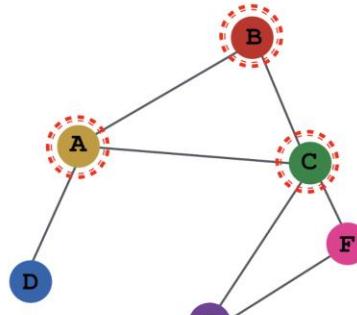
# Model Design: Overview

(1) Define a neighborhood aggregation function



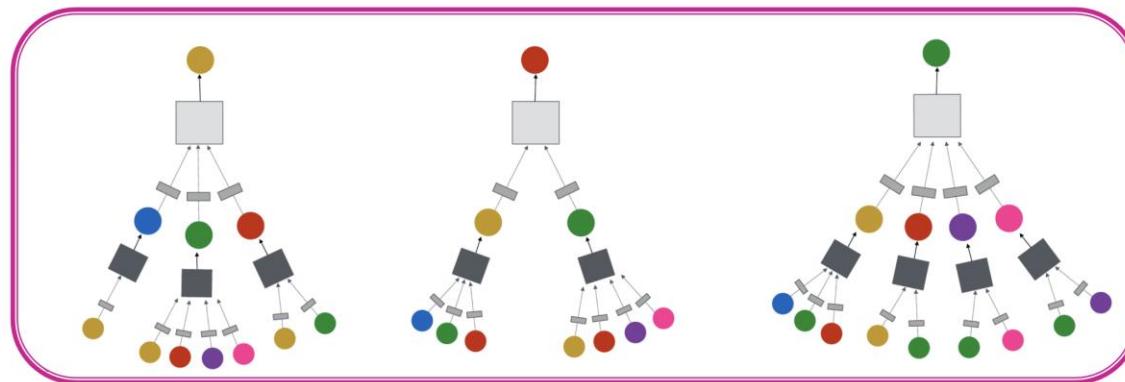
(2) Define a loss function on the embeddings

# Model Design: Overview

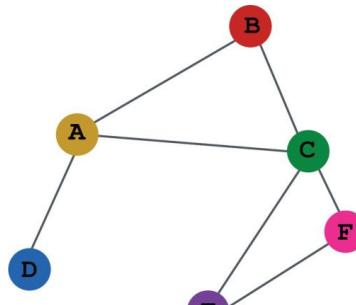


INPUT GRAPH

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



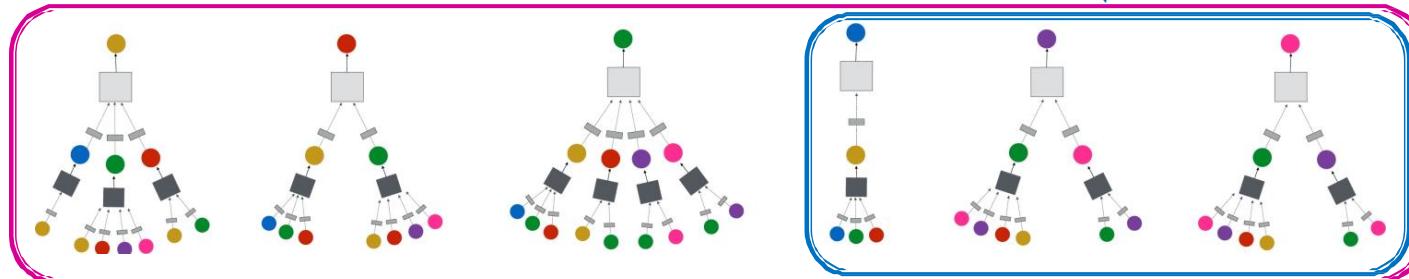
# Model Design: Overview



INPUT GRAPH

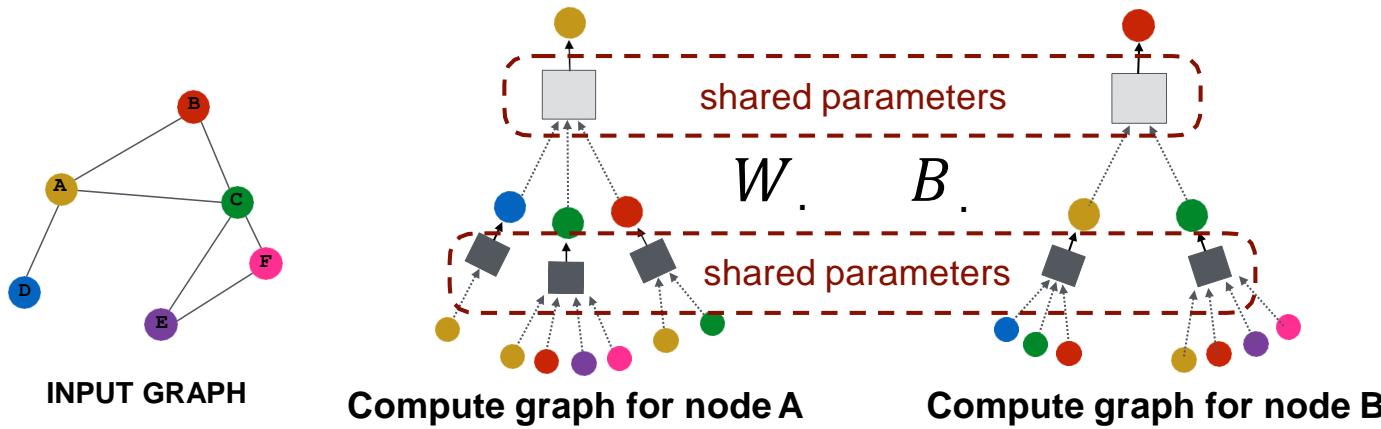
(4) Generate embeddings  
for nodes as needed

Even for nodes we never  
trained on!

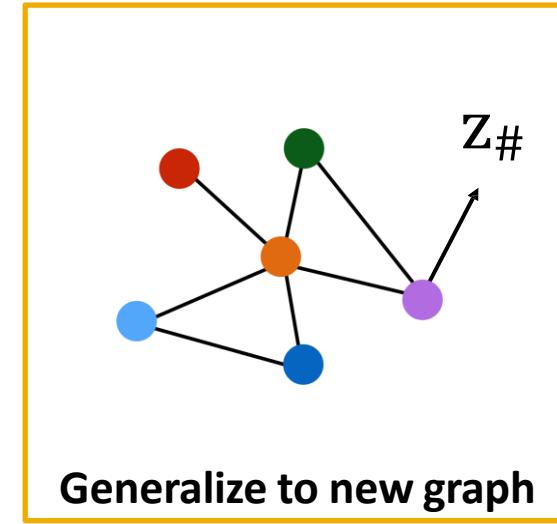
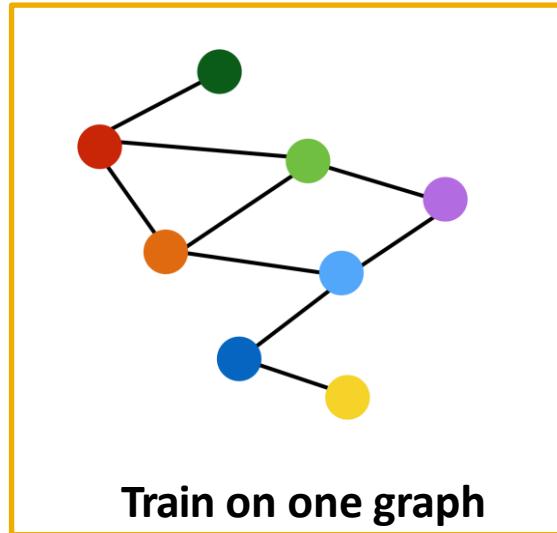


# Inductive Capability

- 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!**



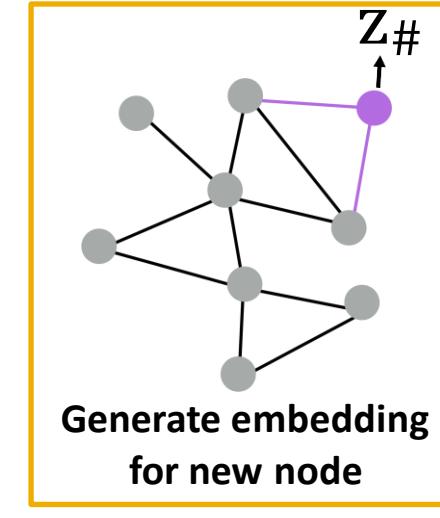
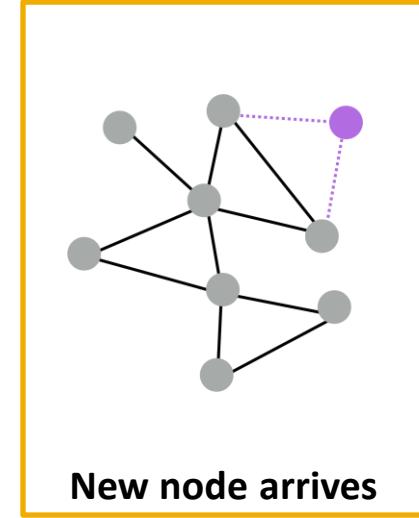
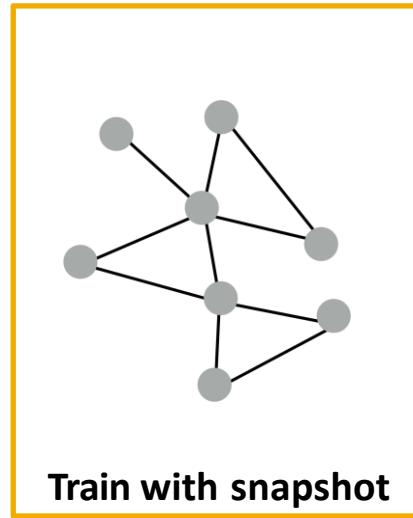
# Inductive Capability: New Graphs



Inductive node embedding → 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

# Inductive Capability: New Nodes



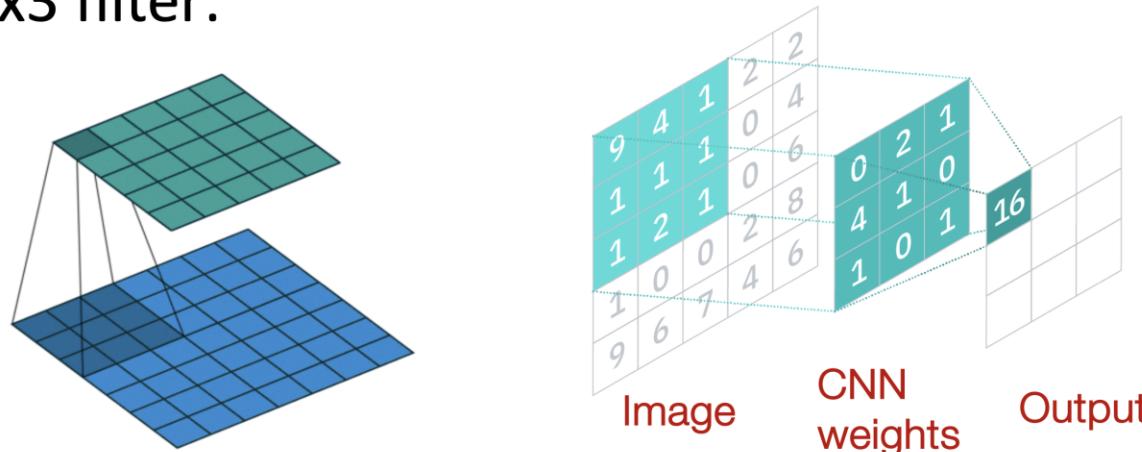
- Many application settings constantly encounter previously unseen nodes:
  - E.g., Reddit, YouTube, Google Scholar
- Need to generate new embeddings “on the fly”

# Architecture Comparison

- How do GNNs compare to prominent architectures such as Convolutional Neural Nets?

# Convolutional Neural Network

Convolutional neural network (CNN) layer with 3x3 filter:

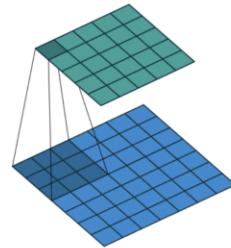


$$\text{CNN formulation: } h_v^{(l+1)} = \sigma(\sum_{u \in N(v) \cup \{v\}} W_l^u h_u^{(l)}), \quad \forall l \in \{0, \dots, L-1\}$$

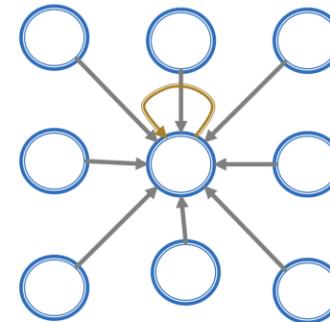
$N(v)$  represents the 8 neighbor pixels of  $v$ .

# GNN vs. CNN

Convolutional neural network (CNN) layer with  
3x3 filter:



Image

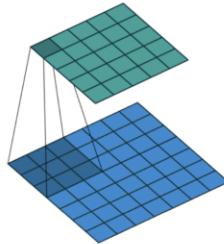


Graph

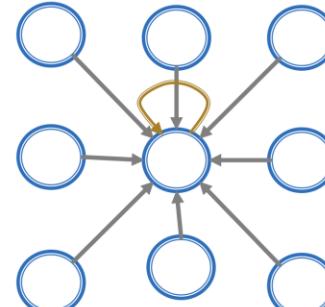
- GNN formulation:  $h_v^{(l+1)} = \sigma(\mathbf{W}_l \sum_{u \in N(v)} \frac{h_u^{(l)}}{|N(v)|} + B_l h_v^{(l)})$ ,  $\forall l \in \{0, \dots, L-1\}$
- CNN formulation: (previous slide)  $h_v^{(l+1)} = \sigma(\sum_{u \in N(v) \cup \{v\}} W_l^u h_u^{(l)})$ ,  $\forall l \in \{0, \dots, L-1\}$   
if we rewrite:  $h_v^{(l+1)} = \sigma(\sum_{u \in N(v)} \mathbf{W}_l^u h_u^{(l)} + B_l h_v^{(l)})$ ,  $\forall l \in \{0, \dots, L-1\}$

# GNN vs. CNN

Convolutional neural network (CNN) layer with  
3x3 filter:



Image



Graph

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

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

**Key difference:** We can learn different  $W_l^u$  for different “neighbor”  $u$  for pixel  $v$  on the image. The reason is we can pick an order for the 9 neighbors using **relative position** to the center pixel:  $\{(-1, -1), (-1, 0), (-1, 1), \dots, (1, 1)\}$

# GNN vs. CNN

Convolutional neural network (CNN) layer with 3x3 filter:

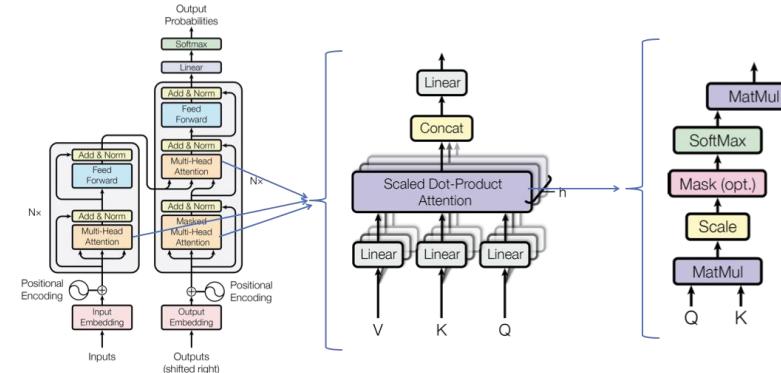


- CNN can be seen as a special GNN with fixed neighbor size and ordering:
  - The size of the filter is pre-defined for a CNN.
  - The advantage of GNN is it processes arbitrary graphs with different degrees for each node.
- CNN is not permutation invariant/equivariant.
  - Switching the order of pixels leads to different outputs.

**Key difference:** We can learn different  $W_v^u$  for different “neighbor”  $u$  for pixel  $v$  on the image. The reason is we can pick an order for the 9 neighbors using **relative position** to the center pixel:  $\{(-1,-1), (-1,0), (-1, 1), \dots, (1, 1)\}$

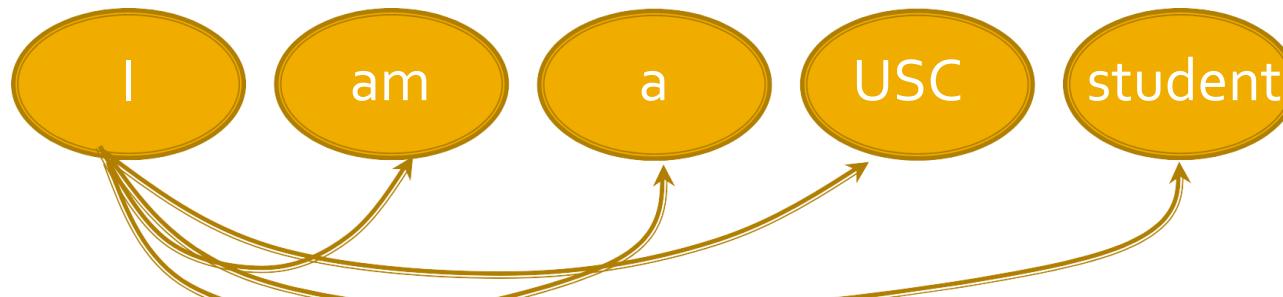
# Transformer

Transformer is one of the most popular architectures that achieves great performance in many sequence modeling tasks.



## Key component: self-attention

- Every token/word attends to all the other tokens/words via matrix calculation.

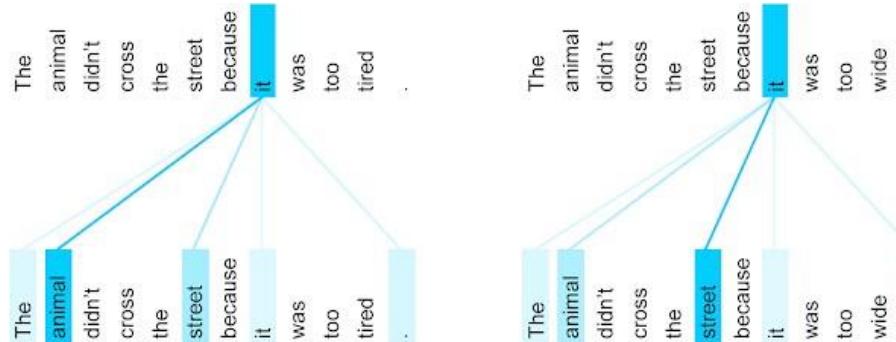


# Transformer

## A general definition of attention:

Given a set of vector *values*, and a vector *query*, **attention** is a technique to compute a weighted sum of the values, dependent on the query.

Each token/word has a **value vector** and a **query vector**. The value vector can be seen as the representation of the token/word. We use the query vector to calculate the attention score (weights in the weighted sum).



# GNN vs. Transformer

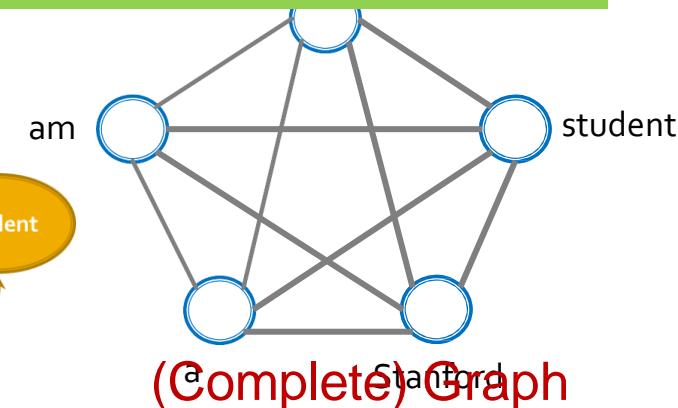
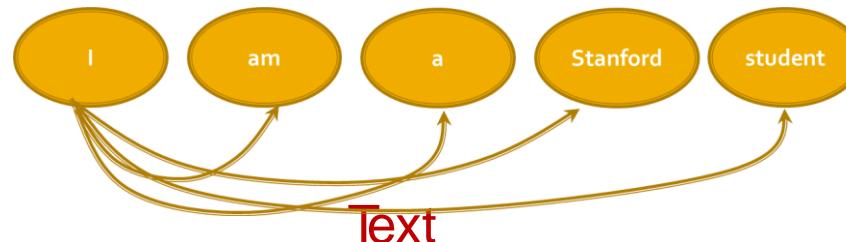
Key



Transformer layer can be seen as a special GNN that runs on a fully-connected “word” graph!

Since

words, **the computation graph** of a transformer layer is identical to that of a GNN on the **fully-connected “word” graph**.



# Summary

- In this lecture, we introduced
  - Idea for Deep Learning for Graphs
    - Multiple layers of embedding transformation
    - At every layer, use the embedding at previous layer as the input
    - Aggregation of neighbors and self-embeddings
  - Graph Convolutional Network
    - Mean aggregation (can be expressed in matrix form)
  - GNN is a general architecture
    - CNN can be viewed as a special GNN