# Deep Learning (1470)
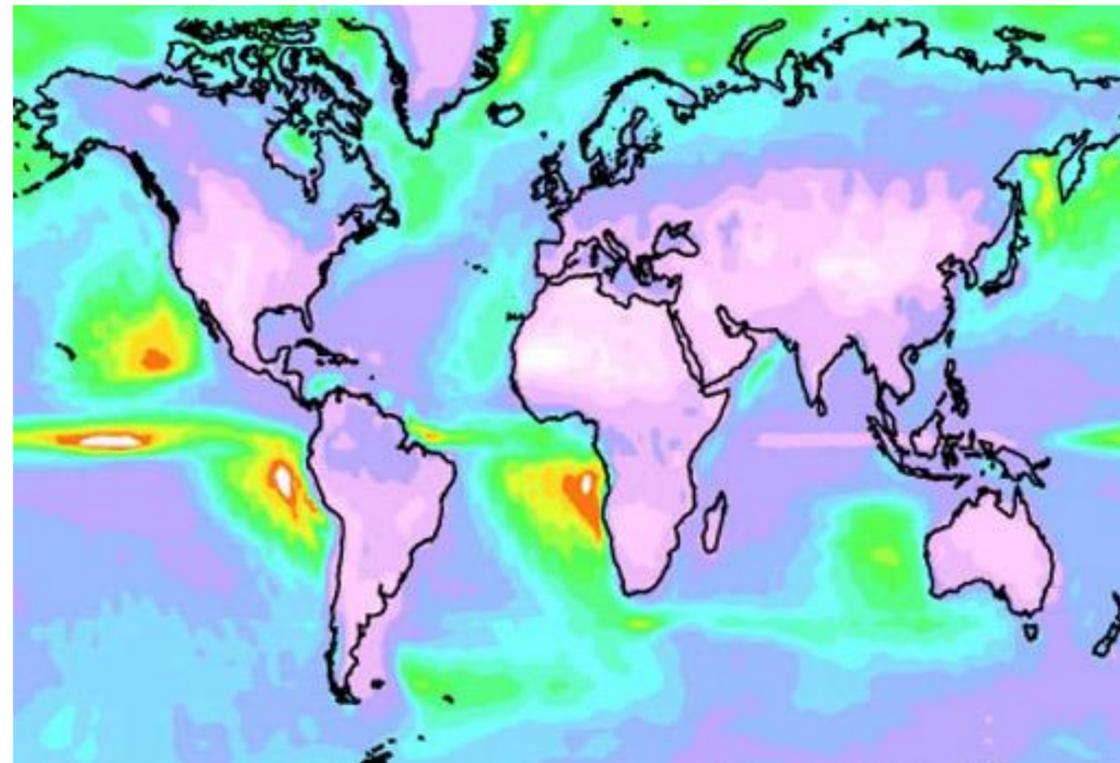
**Randall Balestriero**

**Class 18: Geometric Deep Learning**

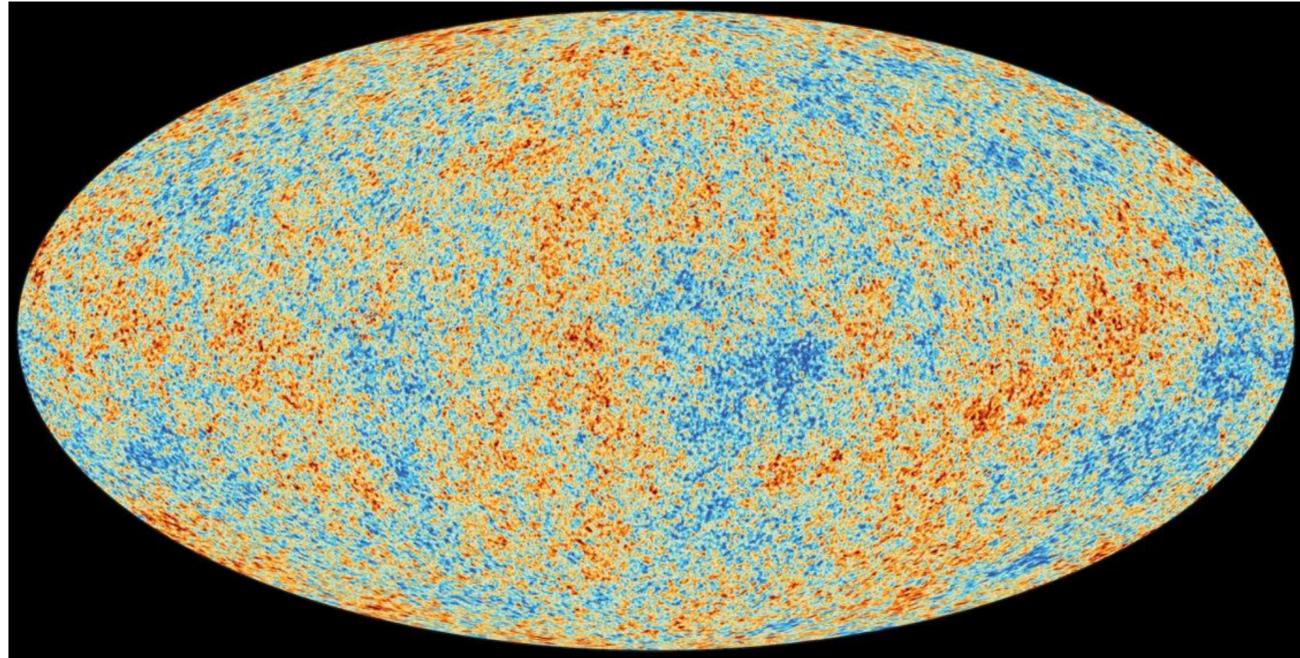# Recap!

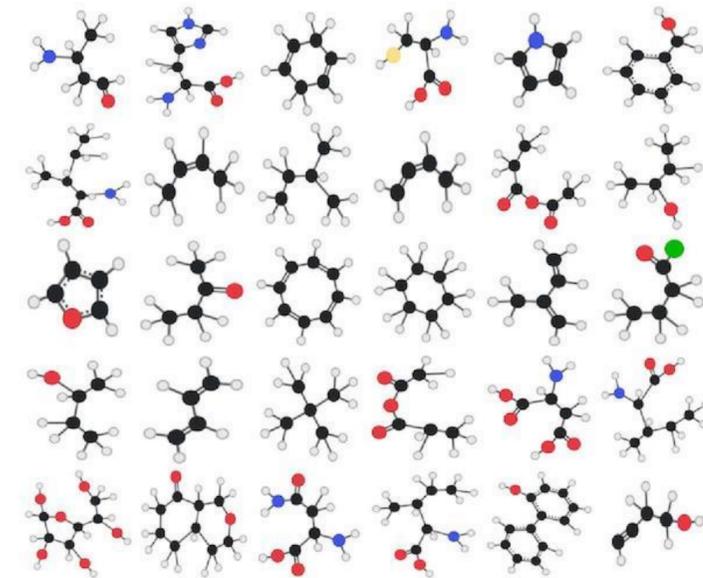Goal: Take as input data from satellites, predict carbon dioxide

Idea 1: Take all of the data as input, make it look like an image, use CNNs to learn to output carbon monoxide amounts
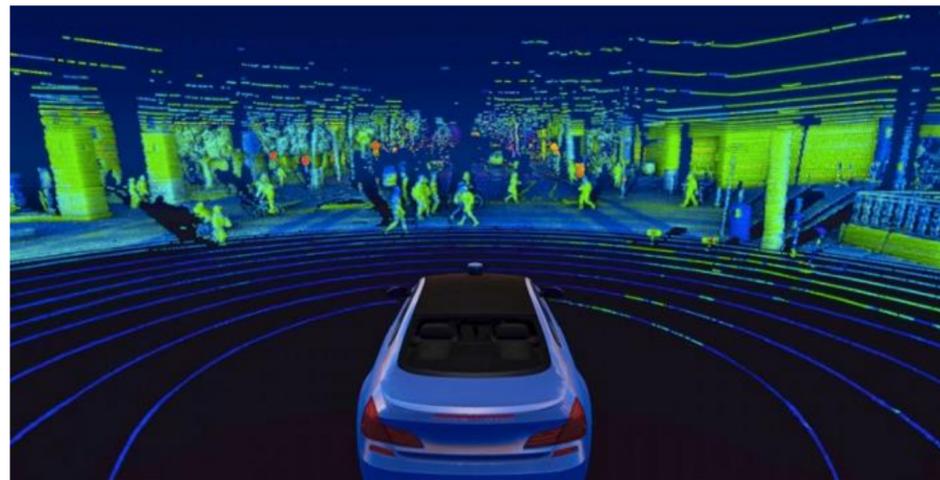
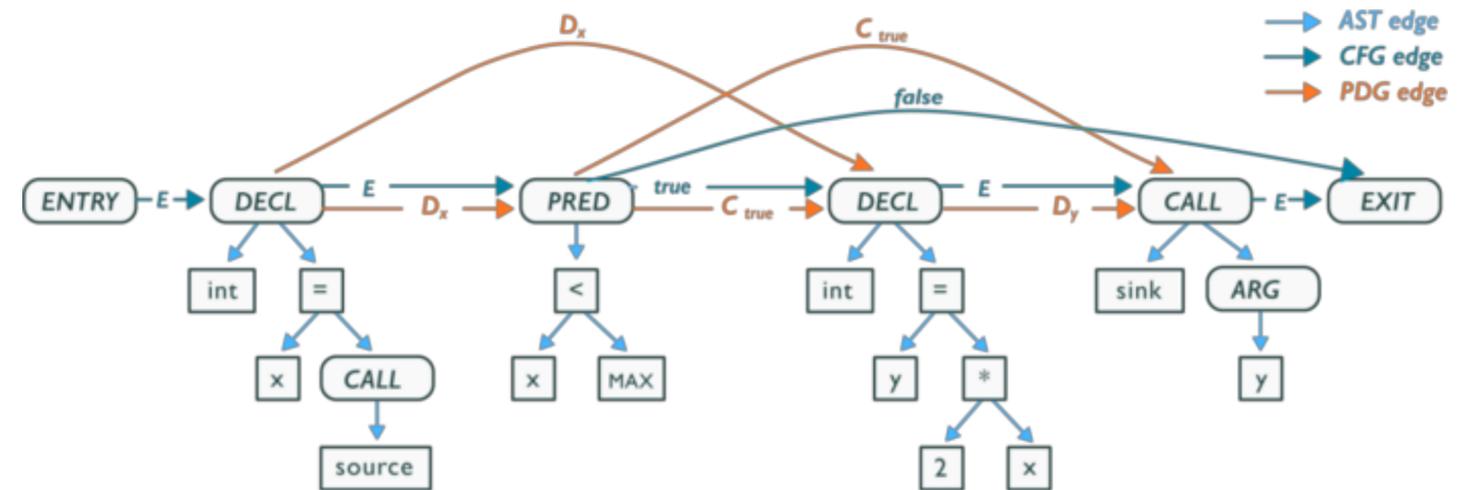# Other Types of Structured Data Types



Microwave Background Radiation
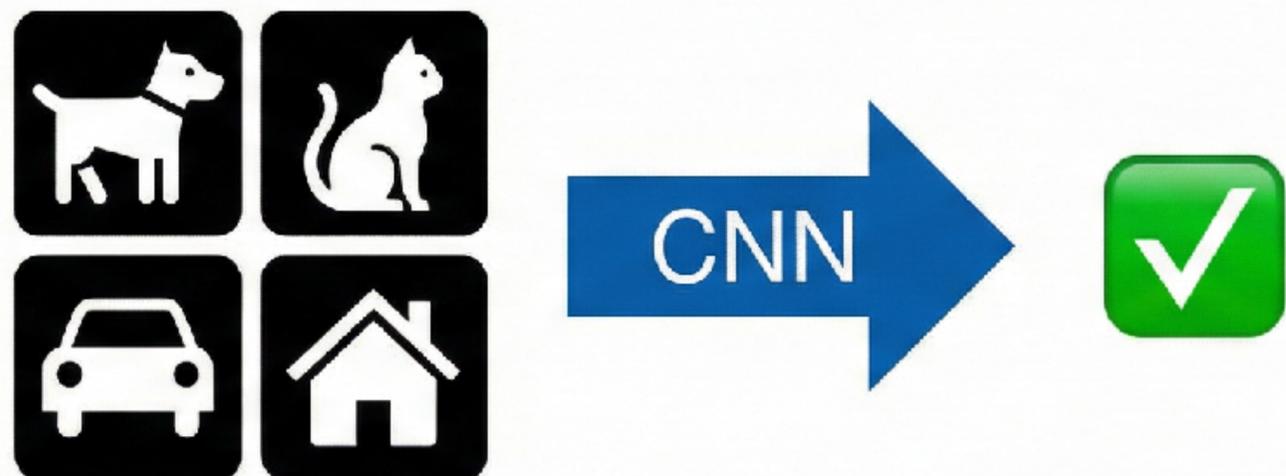


Molecule Data



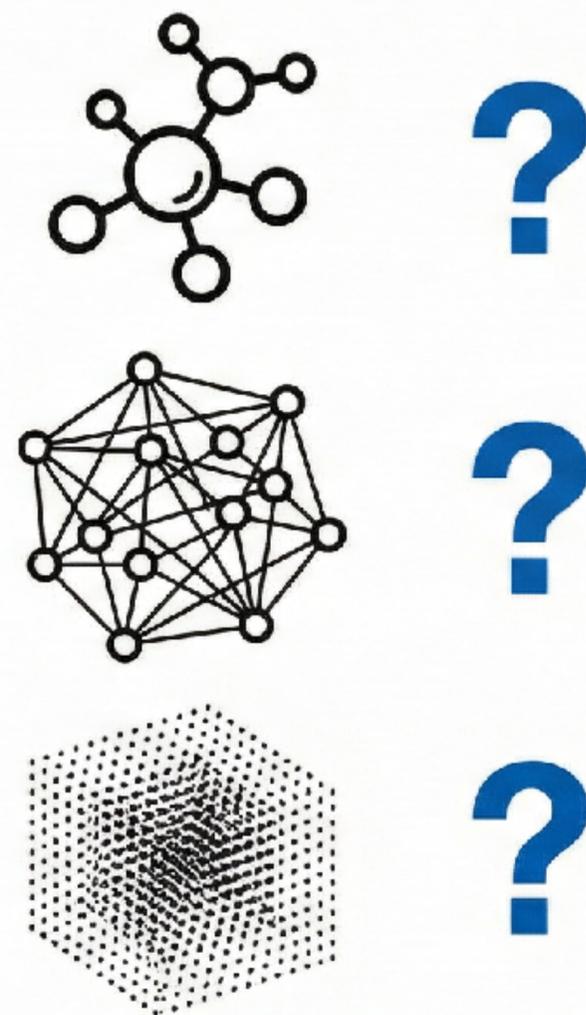Point Clouds (from LIDAR)



Code Graphs
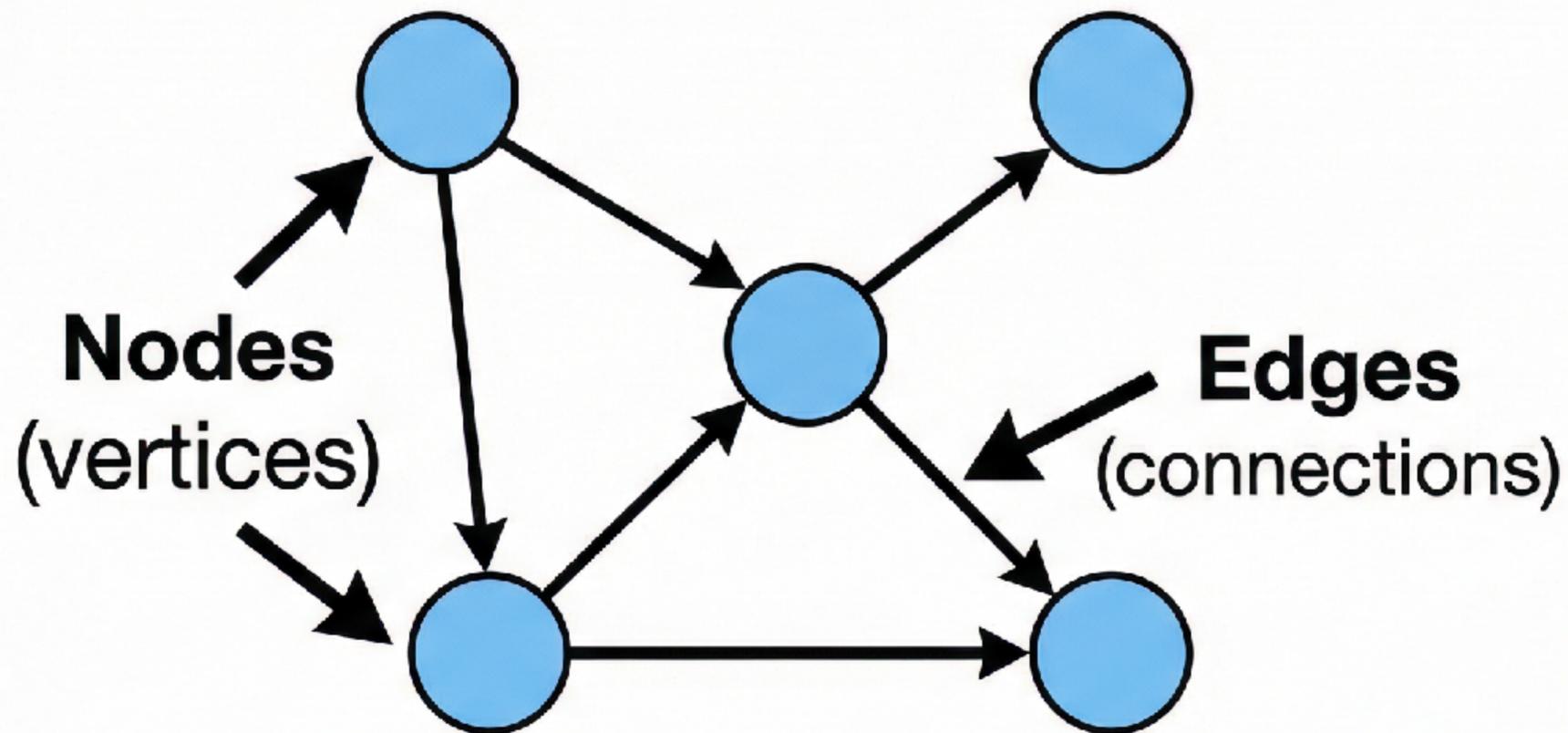
# THE BIG PICTURE

Works

Challenge

CNN ✅

Data with structure needs special networks

# Geometric Deep Learning

- We turned to convolutions for image data to give our networks "spatial reasoning"

- By Explicitly modelling the relationship of our data, we can achieve better results

- Geometric Deep Learning is the subfield of DL dedicated to learning representations of *structured* data.

# GRAPHS 101: A QUICK INTRO

## What is a Graph?



**Nodes** (vertices)

**Edges** (connections)
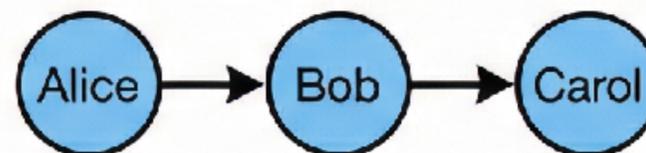
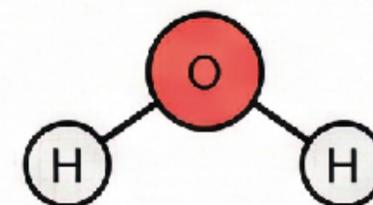Like a map of relationships!

## Examples

Small social network:



Alice - Bob - Carol

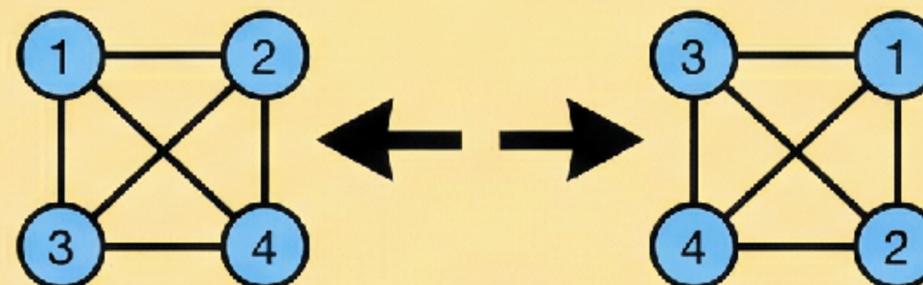Small molecule:



H - O - H

Small citation:



Paper A cites Paper B

### Key Insight
Graphs have NO fixed ordering of nodes!



Both represent the SAME graph
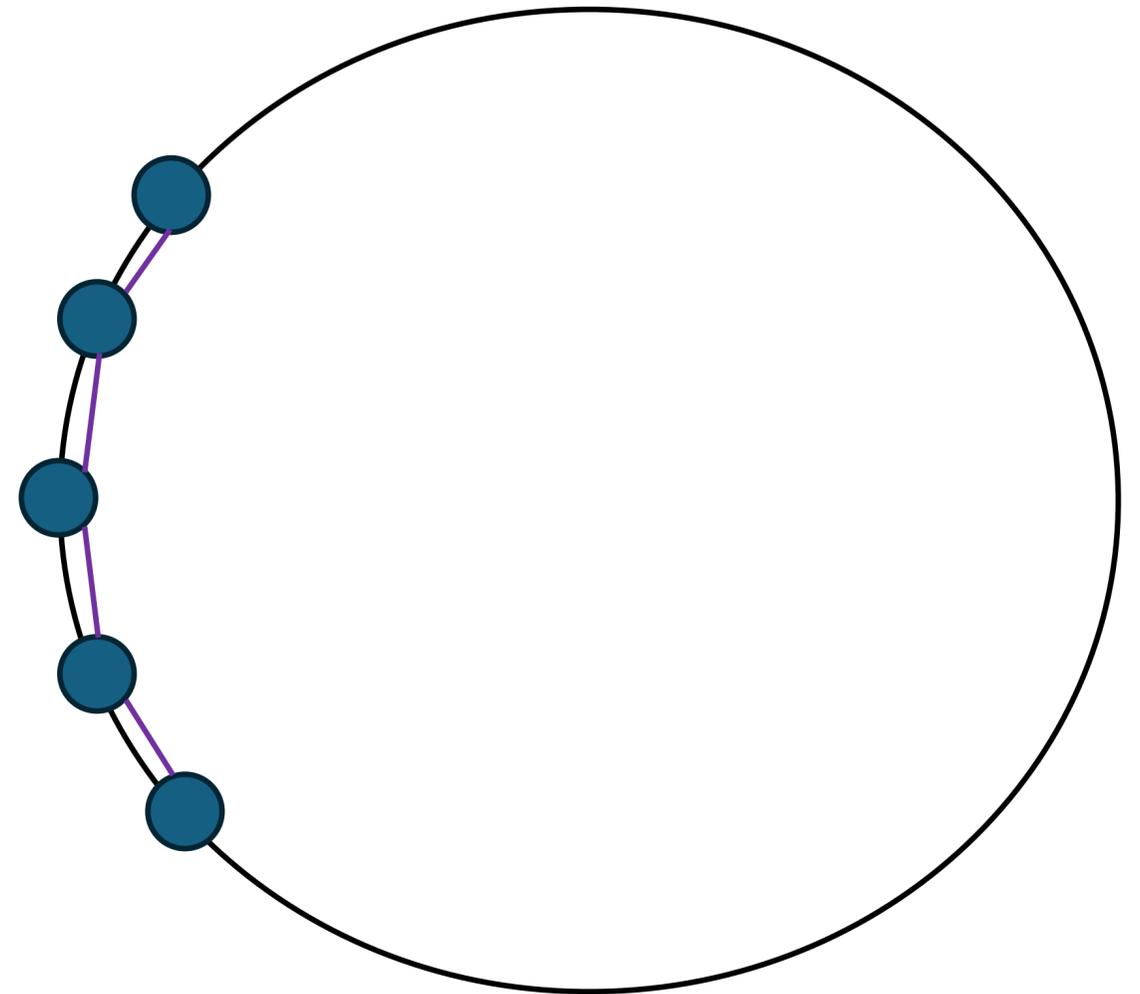
This is **PERMUTATION INVARIANCE**

# Extending Convolutions to Manifolds

A circle is a manifold in 2D
(like a sphere or taurus in 3D)

For every point on the manifold,
connect it to "close" points

What data structure does this lead to?

# HOW DO GRAPH NEURAL NETWORKS WORK?



**STEP 1 - START**

FEATURE VECTOR

Each node has features (like a profile)

**STEP 2 - TALK TO NEIGHBORS**

Each node collects info from neighbors

💬 Like asking friends for advice!

**STEP 3 - UPDATE**

Combine and update your own features

**STEP 4 - REPEAT**

Stack layers = see farther in graph

After K layers, each node knows about nodes K hops away!

# Node Prediction

Each node *v* has a learned representation $z_v$, we can learn a fully-connected layer to go from features $z_v$ to output

$$f(z_v) = \sigma(W z_v)$$

# Link Prediction

Learn a function that takes two nodes as input and predicts the presence (or absence) of an edge

$$f(z_v, z_u) = \sigma(W(z_v \odot z_u))$$

Element-wise product

# GCN: THE CORE EQUATION

$$H' = \sigma\left(D^{-1/2}\,A\,D^{-1/2}\,H\,W\right)$$

## A = Adjacency Matrix
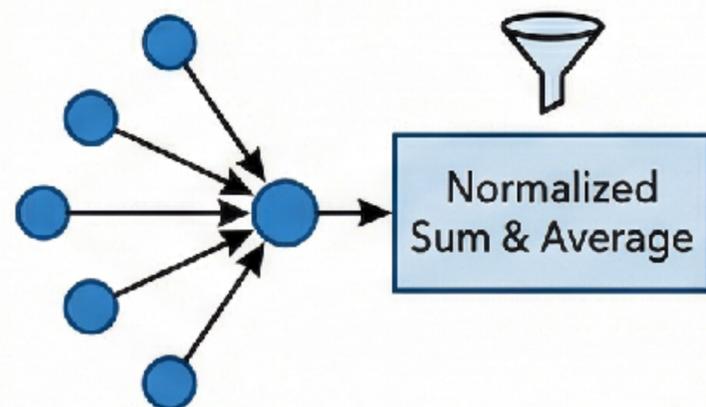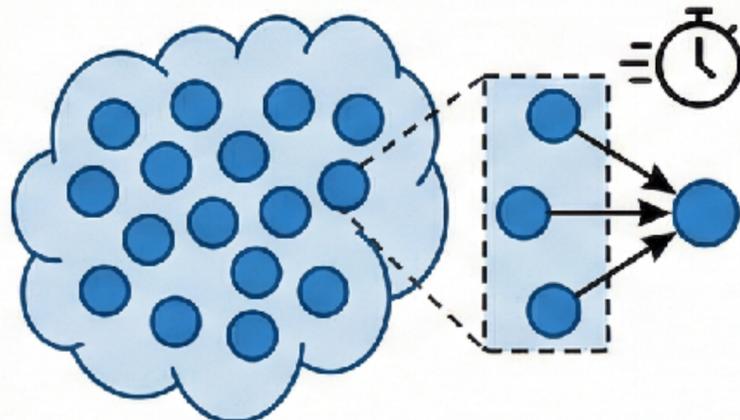
| 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 |

Who connects to whom

## D = Degree Matrix

| 2 | 0 | 0 | 0 |
| 0 | 3 | 0 | 0 |
| 0 | 0 | 2 | 0 |
| 0 | 0 | 0 | 2 |

How many neighbors
each node has

## H, W = Features and Weights

**H (Features)**

N x F

| 4 | 1 | 3 |
| 1 | 2 | 2 |
| ⋮ | ⋮ | ⋮ |
| 4 | ⋮ | 3 |

4 x F

**W (Weights)**

F x F'

| 3 | 1 |
| 1 | 2 |
| ⋮ | ⋮ |
| 4 | 3 |

F x F'

Learnable transformation

**Aggregate neighbors  ->  Normalize by degree  ->  Transform  ->  Activate**

# GCN: IMPLEMENTATION DETAILS

## PyTorch Code Structure

```
class GCNLayer:
  - W: weight matrix (F_in, F_out)
  - forward(H, A):
    - A_hat = A + I (add self-loops)
    - D_hat = degree(A_hat)
    - norm = D^(-1/2) A_hat D^(-1/2)
    - return ReLU(norm @ H @ W)
```

## Layer Stacking



2-3 layers typical (avoid over-smoothing)

## Key Implementation Choices

**Activation**
ReLU, LeakyReLU, none on final

**Dropout**
Between layers for regularization

**Residual**
H' = GCN(H) + H helps with depth

# GAT: ATTENTION MECHANISM

## The Key Idea
Not all neighbors are equally important!

Learn attention weights for each edge

## The Attention Equation

Step 1: Project features

$$\mathbf{z}_i = \mathbf{W}\mathbf{h}_i$$

(linear transformation)

Step 2: Compute attention scores

$$e_{ij} = \text{LeakyReLU}\left(\mathbf{a}^T\left[\mathbf{z}_i \,\|\, \mathbf{z}_j\right]\right)$$

Concatenate and score with vector **a**

Step 3: Normalize with softmax

$$\alpha_{ij} = \text{softmax}_j(e_{ij}) = \frac{\exp(e_{ij})}{\sum_k \exp(e_{ik})}$$

Weights sum to 1 over neighbors

Final aggregation

$$\mathbf{h}_i' = \sigma\left(\sum_j \alpha_{ij}\mathbf{W}\mathbf{h}_j\right)$$

Weighted sum of neighbor features

# GAT: MULTI-HEAD ATTENTION

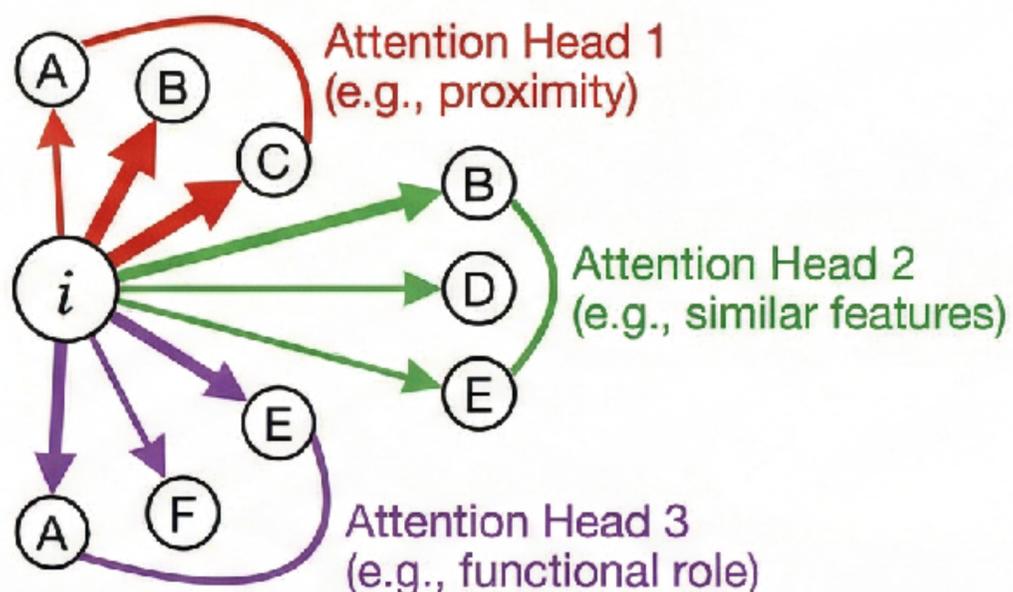## Why Multiple Heads?

Different heads learn different relationship types!



Attention Head 1
(e.g., proximity)

Attention Head 2
(e.g., similar features)

Attention Head 3
(e.g., functional role)

Like having multiple perspectives

## Multi-Head Formulation

$$\mathbf{h}_i' = \text{CONCAT}(\text{head}_1, \text{head}_2, ..., \text{head}_K) \text{ or } \text{MEAN}(\text{heads})$$

$$\text{Each head}_k = \sigma\left(\sum_j \alpha_{ij}^k \mathbf{W}^k \mathbf{h}_j\right)$$

$K$ independent attention mechanisms
Concatenate (middle layers) or Average (final layer)
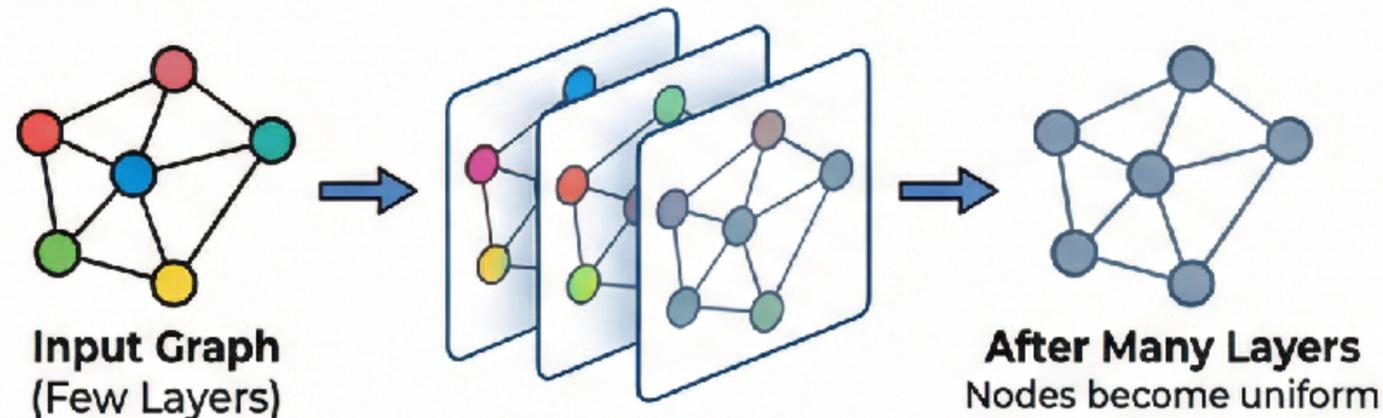
## Implementation Summary

### Learnable Parameters

- $W$: feature projection ($F$ x $F'$)
- $a$: attention vector ($2F'$ x 1)
- Per head: $W^k$ and $a^k$

### Architecture Choices

- Heads $K$: typically 8
- Hidden dim $F'$: 8 per head
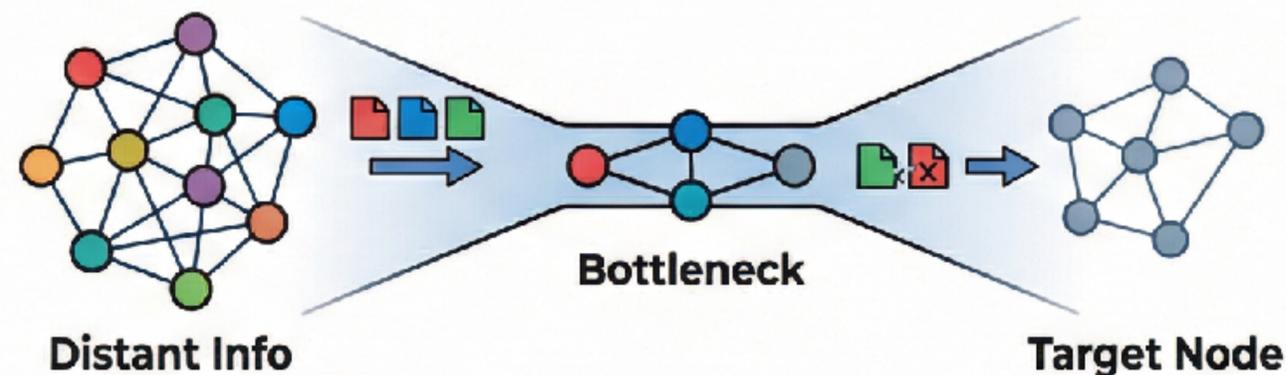- Dropout on attention weights
- Add self-loops for self-attention
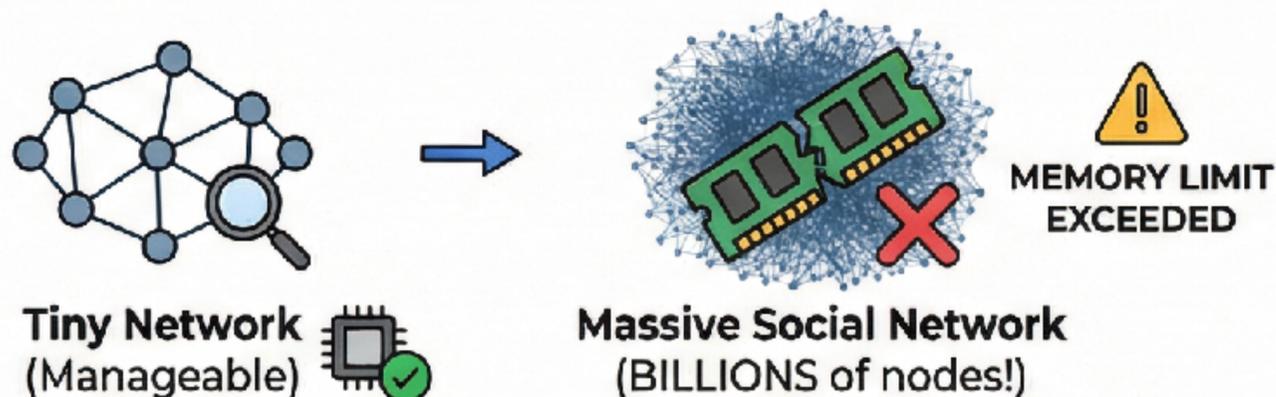
# CHALLENGES IN GEOMETRIC DL

## OVER-SMOOTHING

**Input Graph** (Few Layers) → **After Many Layers** Nodes become uniform

- Too many layers = all nodes look the same!
- Like a game of telephone - info gets blurry

## OVER-SQUASHING

**Distant Info** → **Bottleneck** → **Target Node**
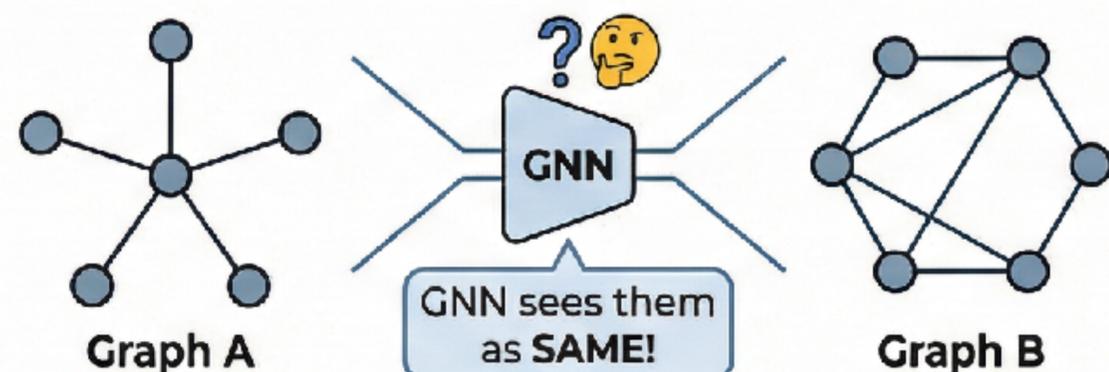
- Hard to send info through narrow graph paths
- Distant nodes can't communicate well

## SCALABILITY

**Tiny Network** (Manageable) → **Massive Social Network** (BILLIONS of nodes!) **MEMORY LIMIT EXCEEDED**

- Social networks have BILLIONS of nodes
- Can't fit in memory!

## EXPRESSIVITY

**Graph A** — GNN — **Graph B**

GNN sees them as SAME!

- Some graphs GNNs cannot tell apart!
- Fundamental theoretical limits

Active research areas - lots of room for new ideas!

See you on Friday!