



# Deep Learning on Graphs

- Introduction of Graph, Graph Machine Learning, Graph Neural Network, & GraphStorm

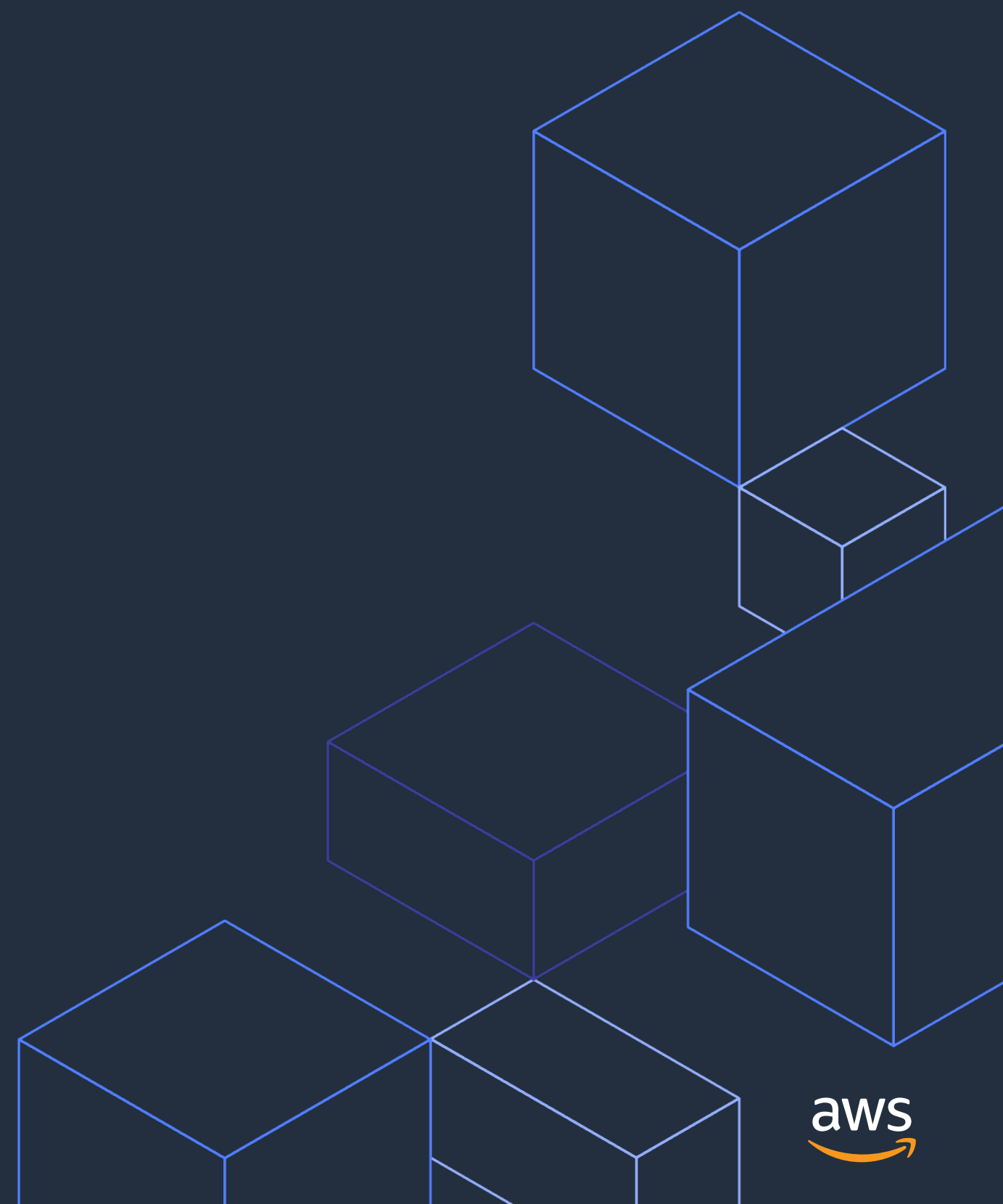
**Jian Zhang**

**AWS AI Research & Education**

# Agenda

- Graph
- Graph Machine Learning (GML)
- Graph Neural Network (GNN)
- Using GNN in Real World
- GraphStorm Framework
- Break
- Hands on GraphStorm

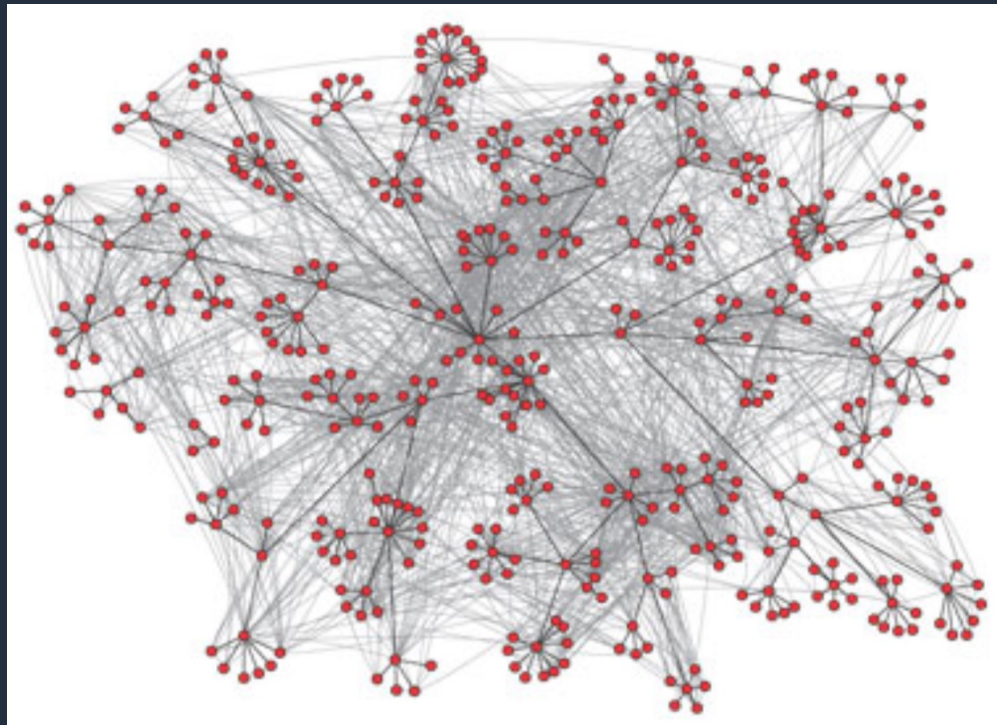
# Graph





# Concepts and terms

## Graph vs Image



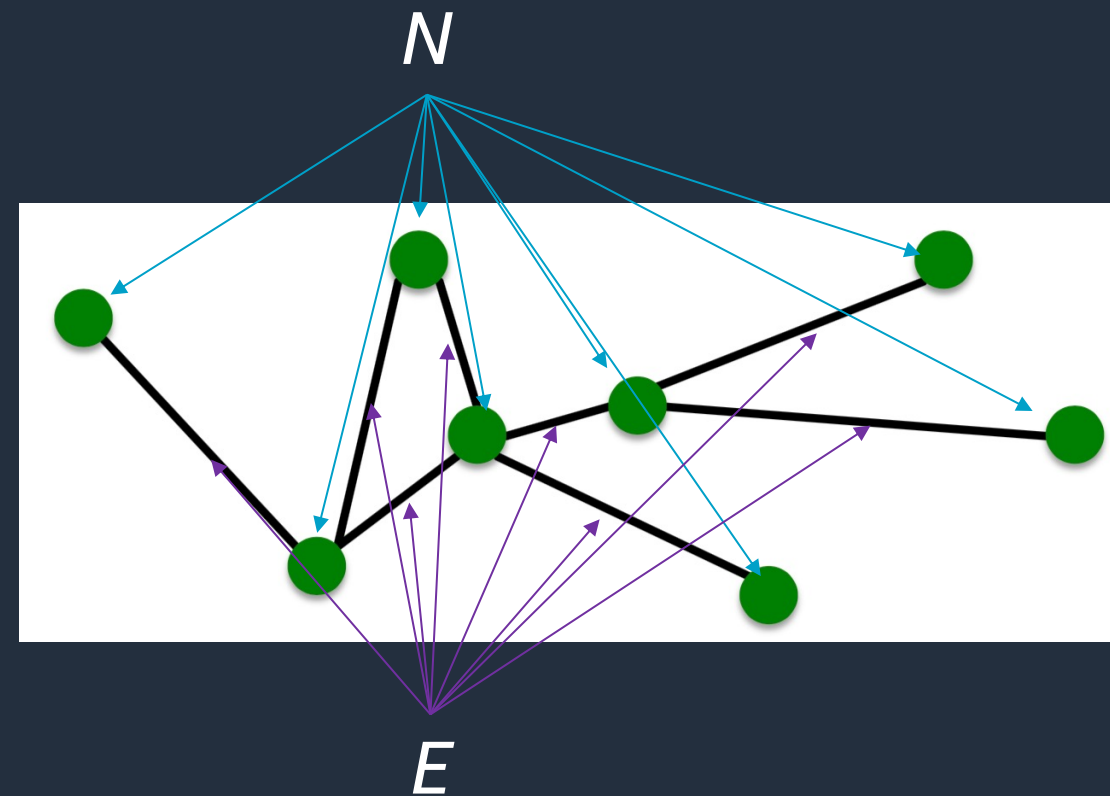
“图”



# Concepts and terms

Node vs Vertex

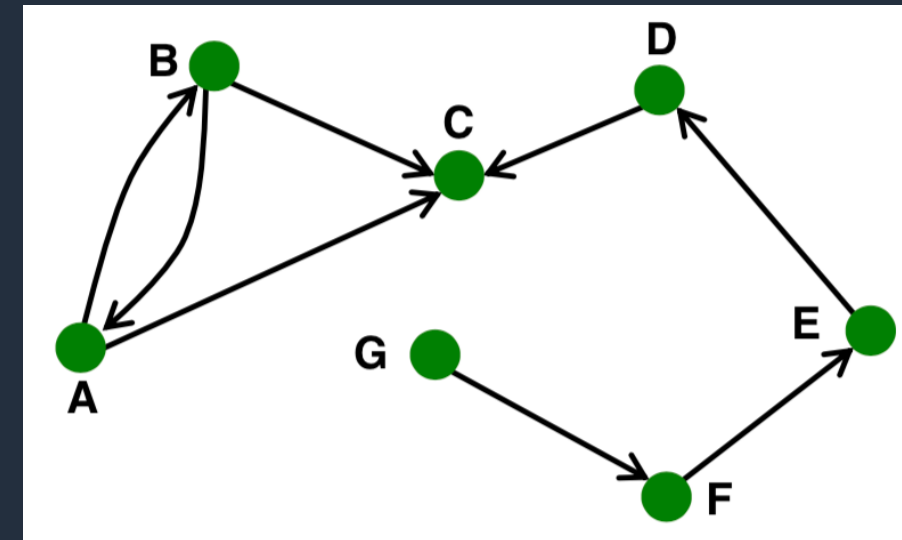
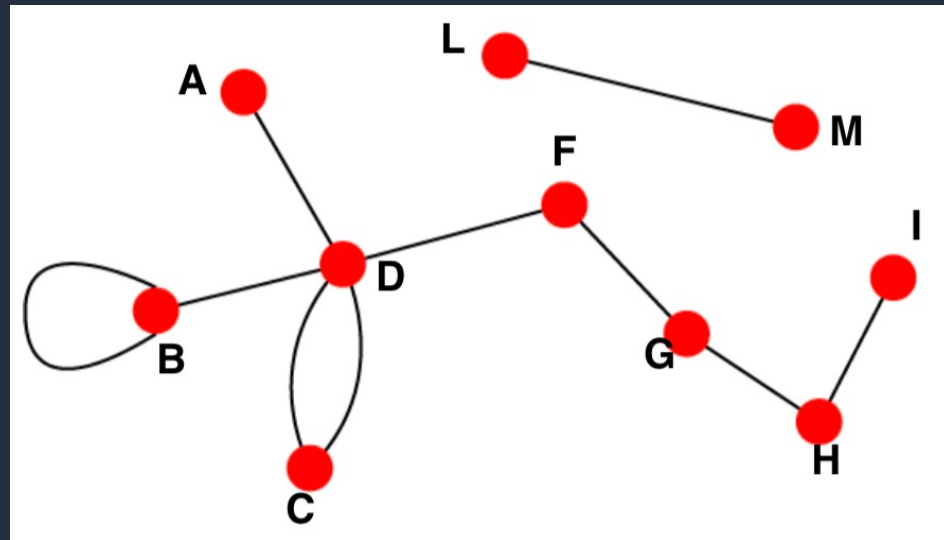
Edge vs Link



→  $G(N, E)$

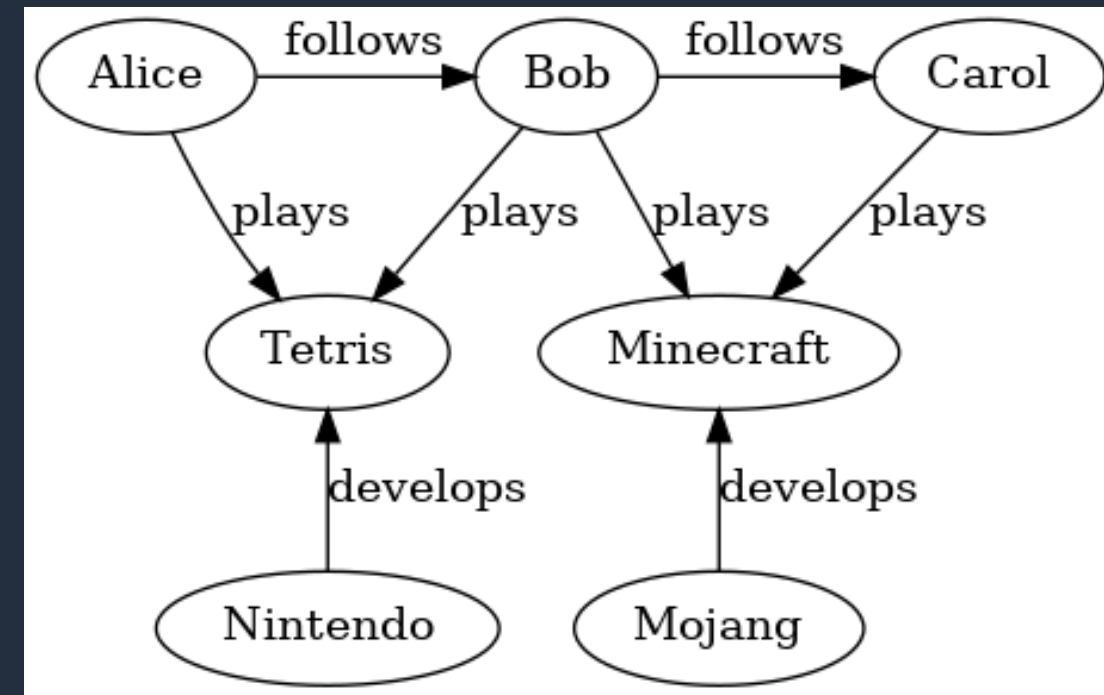
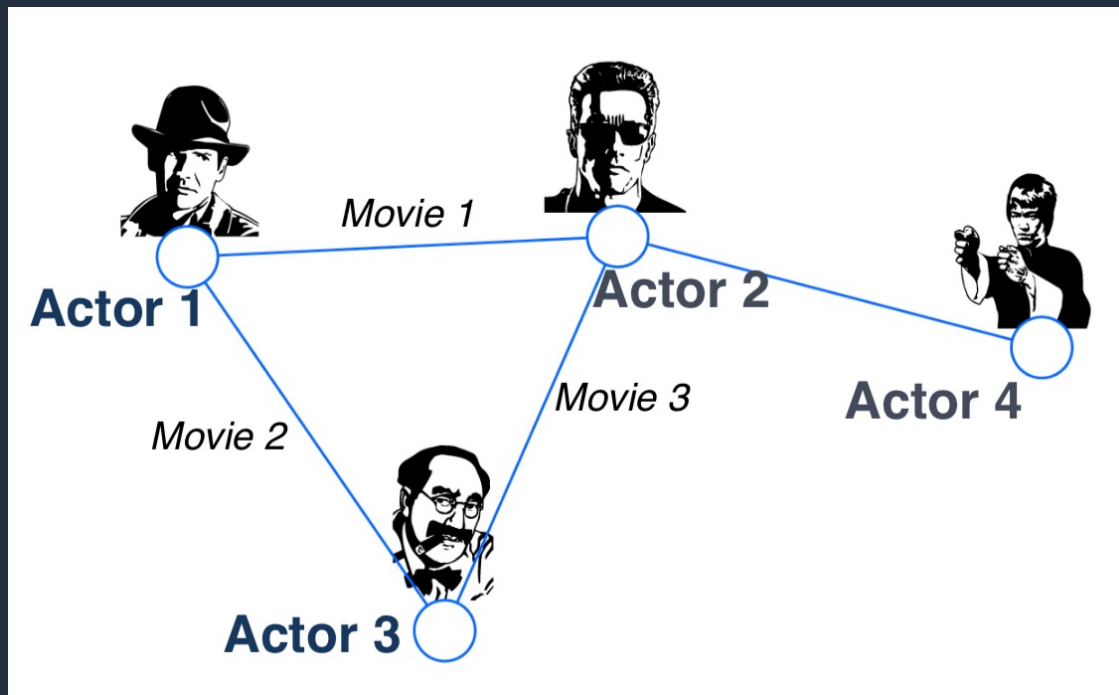
# Concepts and terms

## Undirected vs Directed



# Concepts and terms

## Homogeneous vs Heterogeneous





# Concepts and terms

## Adjacency Matrix vs Edge List



$A_{ij} = 1$  if there is a link from node  $i$  to node  $j$

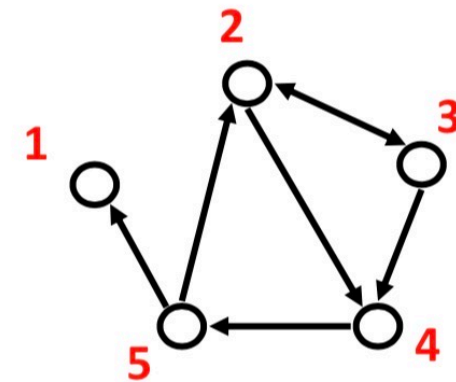
$A_{ij} = 0$  otherwise

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

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

Represent graph as a set of edges:

- (2, 3)
- (2, 4)
- (3, 2)
- (3, 4)
- (4, 5)
- (5, 2)
- (5, 1)



# Concepts and terms

## Features:

Graph Edge List	
src	dst
U1	P1
U1	P2
U1	U4
...	...
U3	P11
U5	P11
U5	P12

## Nodes

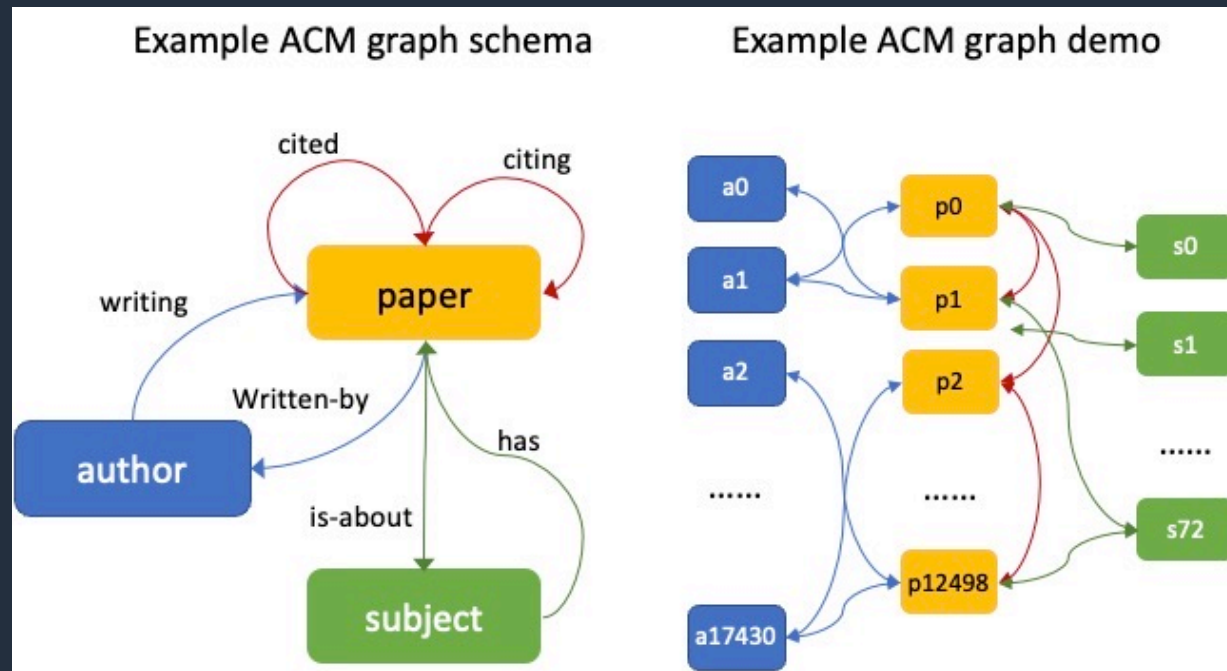
Node Features				
Node	Label	nf1	nf2	...
U1	0	10	0.5	
P1				
...		...	...	
U5	1	50	0.9	
P12				

## Edges

Edge Features				
src	dst	ef1	ef2	...
U1	P1	0.1	L	
U1	P2	0.2	C	
U1	U4	1.4	C	
...	...	...	...	
U3	P11	...	...	
U5	P11	...	...	
U5	P12	...	...	

# Concepts and terms

## Heterogeneous Graph & Features



### Node Parquets

author node parquet

node_id	feat		
a0	...	...	...
a1	...	...	...
a2	...	256D	...
...	...	...	...
a17430	...	...	...

paper node parquet

node_id	feat			label
p0	...	...	...	11
p1	...	...	...	2
p2	...	256D	...	0
...	...	...	...	...
p12497	...	...	...	11

subject node parquet

node_id	feat		
s0	...	...	...
s1	...	...	...
s2	...	256D	...
...	...	...	...
s71	...	...	...

### Edge Parquets

author, writing, paper edge parquet

source_id	dest_id
a148	p0
a148	p1
...	...
a4653	p11567

paper, citing, paper edge parquet

source_id	dest_id	label
p74	p25	1
p8	p1	1
...	...	...
p963	p6607	1

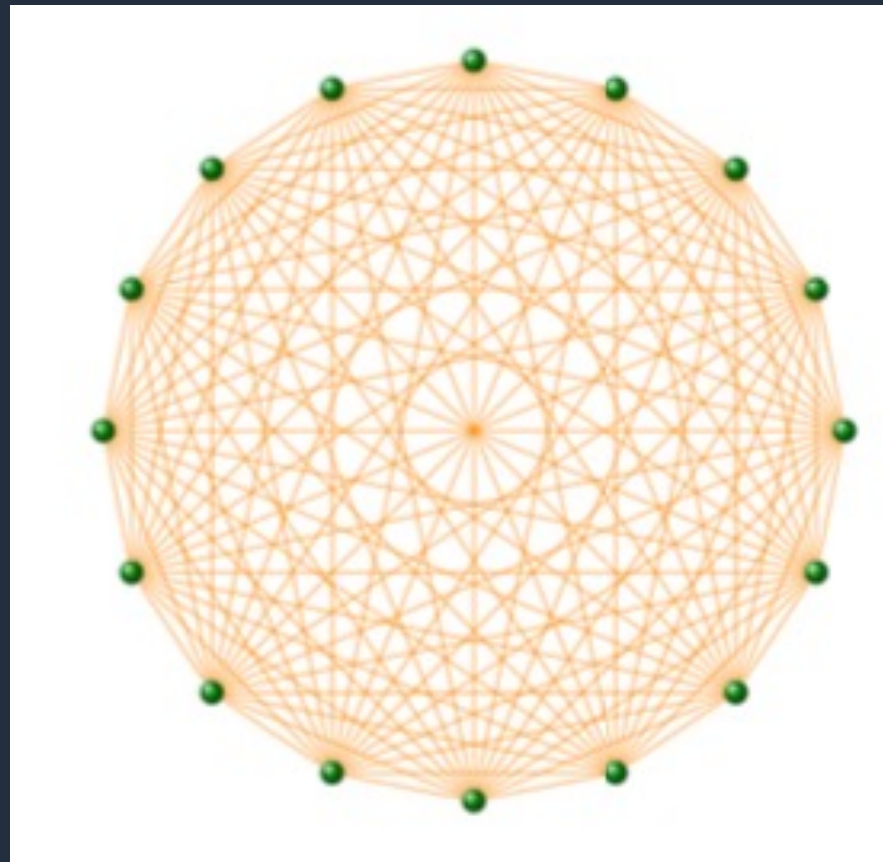
.....

subject, has, paper edge parquet

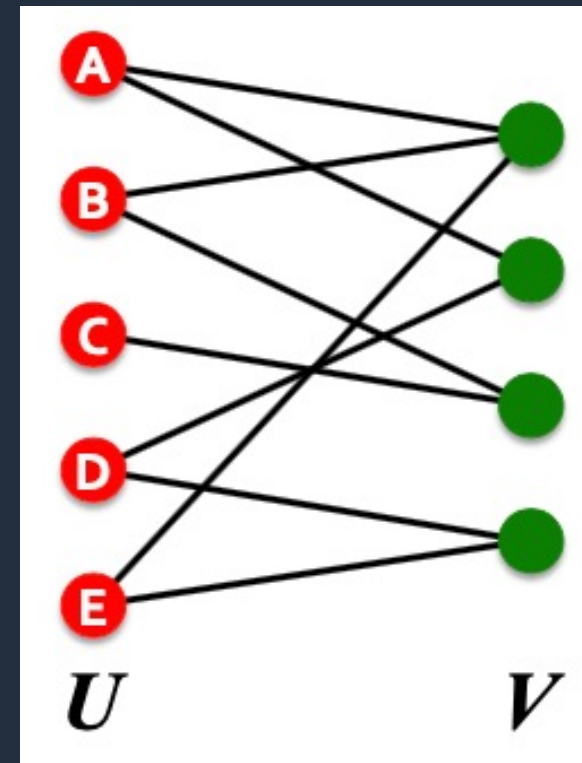
source_id	dest_id
s0	p0
s3	p443
...	...
s56	p11567

# Concepts and terms

## Completed Graph



## Bipartite



# Concepts and terms

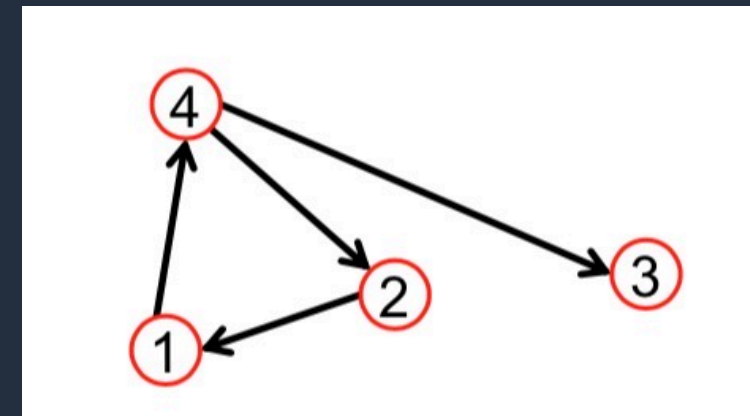
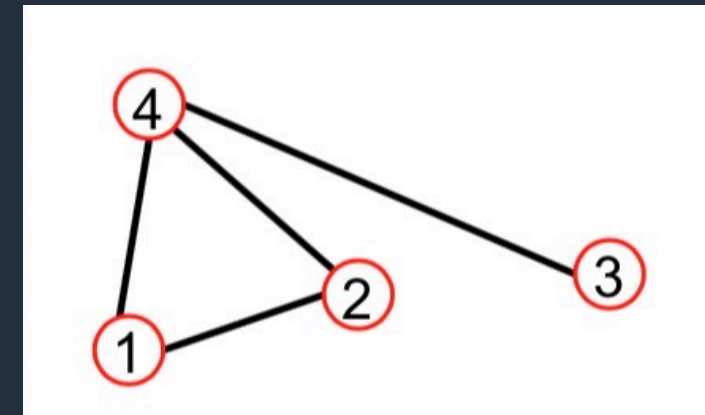
## Node Degree

The number of edges adjacent to a node, e.g.

$$k_4 = 3$$

In directed graphs, a node has an **in-degree** and **out-degree**. And total degree is the sum of in- and out-degrees. e.g.

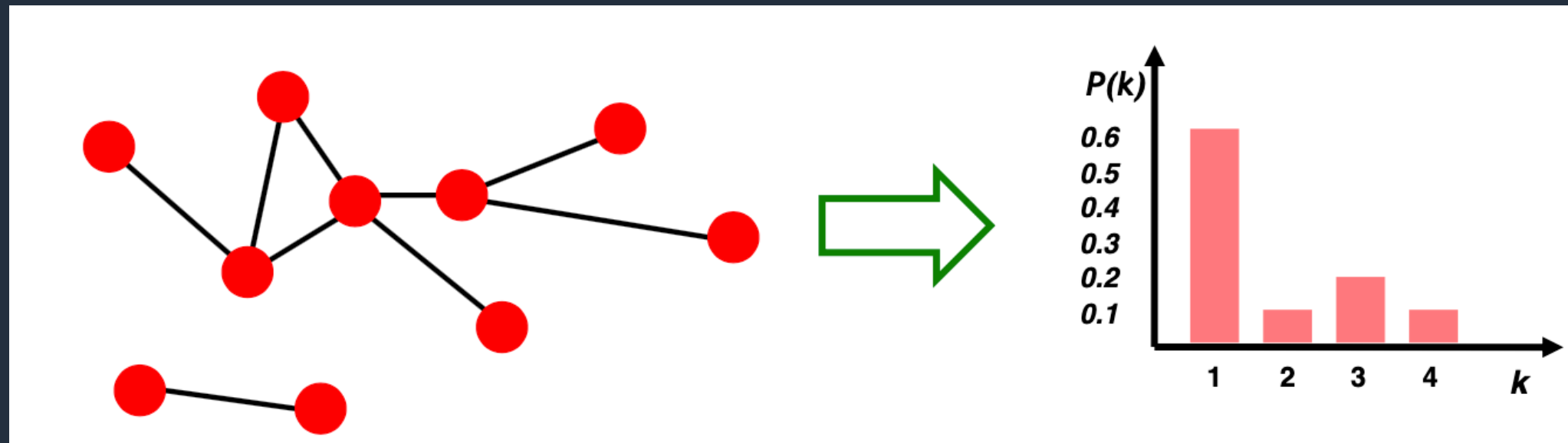
$$k_4^{\text{in}} = 1, k_4^{\text{out}} = 2, \\ k_4^{\text{total}} = 3$$



# Concepts and terms

## Degree Distribution $P(k)$

Probability that a randomly chosen node has degree  $k$ ,  $N_k =$   
*# nodes with degree  $k$* . Normalized  $P(k) = N_k/N$



# Concepts and terms

## Path $P$

A **path** is a sequence of nodes in which each node is linked to the next one.  $P_n = \{i_0, i_1, i_2, \dots, i_n\}$ . Path length  $h = \#$  of edges in a path. In directed graph, paths need to follow the direction of edges.

The **distance** between two nodes is the length of the shortest path.

The **diameter** of a graph is the maximum distance between any pair of nodes.

Average path length of connected graph

The **average path length** for a connected graph (need to be connected).

$$\bar{h} = \frac{1}{2E_{max}} \sum_{i,j \neq i} h_{ij}, \text{ where } h_{ij} \text{ is the distance of node } i \text{ and } j, \text{ and } E_{max} = n(n-1)/2$$

# Concepts and terms

## Cluster Coefficient $C$

A **Cluster Coefficient** is the ratio of the number of edges among neighbors of node  $i$  and the max number of edges between all neighbors:  $C_i = \frac{2e_i}{k_i(k_i - 1)}$ .

**Average Clustering Coefficient** is the mean of all nodes' **Clustering Coefficient**.



# Graph Machine Learning (GML)

# Recall Machine Learning (ML)

$$f(X, \theta) \Rightarrow Y$$

- Logistic Regression
  - Support Vector Machine
  - Multilayer Perceptron
  - Convolutional Neural Network
  - Recurrent Neural Network
  - Transformer
  - ...
- Classification
  - Regression
  - Segmentation
  - Detection
  - Reconstruction
  - ...

The screenshot shows a Jupyter Notebook interface. At the top, the notebook title is 'house\_pricing\_index\_united\_states'. Below the title, there are tabs for 'Data Card', 'Code', and 'Download'. The main content area is divided into two parts. On the left, there is a chat window for ChatGPT. The chat window contains the following text:

**ChatGPT**  
In OpenReview, you typically use LaTeX syntax for mathematical expressions, including superscript and subscript text. Here's how you can do it:

**Superscript:**  
To write text in superscript in OpenReview, you can use the caret symbol '^' followed by curly braces '{ }' to enclose the text you want to be in superscript.

Example:

```
csharp  
This is a^{2} superscript.
```

Output:  
This is a<sup>2</sup> superscript.

**Subscript:**  
To write text in subscript in OpenReview, you can use the underscore symbol '\_' followed by curly braces '{ }' to enclose the text you want to be in subscript.

Example:

```
csharp  
H_{2}O is water.
```

Output:  
H<sub>2</sub>O is water.

In OpenReview, as in LaTeX, you can use the '^' and '\_' without braces '{ }' if the subscript or

On the right side of the chat window, there is a data table with the following content:

Central	DV_ENC
Central	DV_ENC
Central	DV_ENC
Central	DV_ENC
Central	DV_ENC
Central	DV_ENC
Central	DV_ENC

# Graph Machine Learning (GML)

$$f(X, \theta) \Rightarrow Y$$

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

Node Parquets

author node parquet

node_id	feat		
a0	...	...	...
a1	...	...	...
a2	...	256D	...
...	...	...	...
a17430	...	...	...

paper node parquet

node_id	feat			label
p0	...	...	...	11
p1	...	...	...	2
p2	...	256D	...	0
...	...	...	...	...
p12497	...	...	...	11

subject node parquet

node_id	feat		
s0	...	...	...
s1	...	...	...
s2	...	256D	...
...	...	...	...
s71	...	...	...

Edge Parquets

author, writing, paper edge parquet

source_id	dest_id
a148	p0
a148	p1
...	...
a4653	p11567

paper, citing, paper edge parquet

source_id	dest_id	label
p74	p25	1
p8	p1	1
...	...	...
p963	p6607	1

.....

subject, has, paper edge parquet

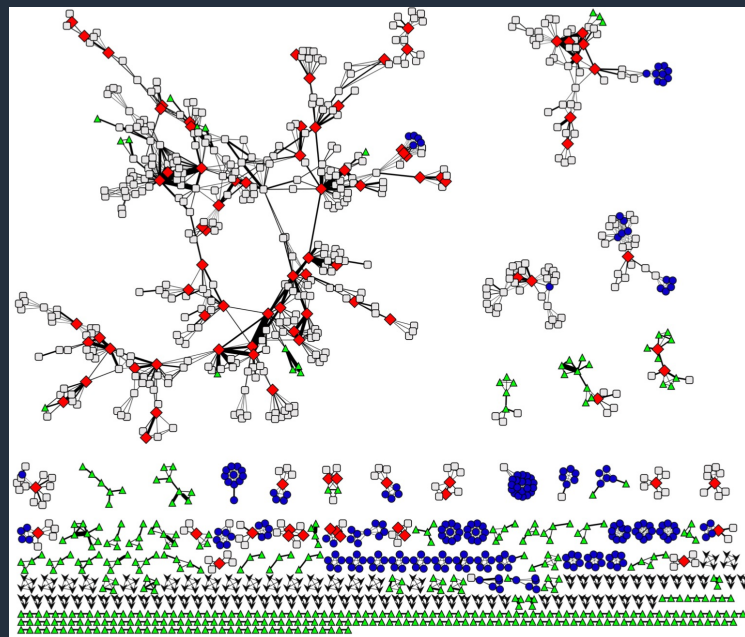
source_id	dest_id
s0	p0
s3	p443
...	...
s56	p11567

- Node
  - Classification, e.g., detecting malicious accounts
  - Regression, e.g., predicting customer rating
- Edge
  - Classification, e.g., detecting suspicious transactions
  - Regression, e.g., predicting when will traffic jam start
  - Link Prediction e.g., recommending friends
- Graph
  - Classification, e.g., predicting if a new compound is toxic
  - Regression e.g., predicting medicine molecular solubility



# GML before

- Generate embeddings by manual feature engineering
- Automatically generate embeddings using unsupervised dimensionality reduction approaches



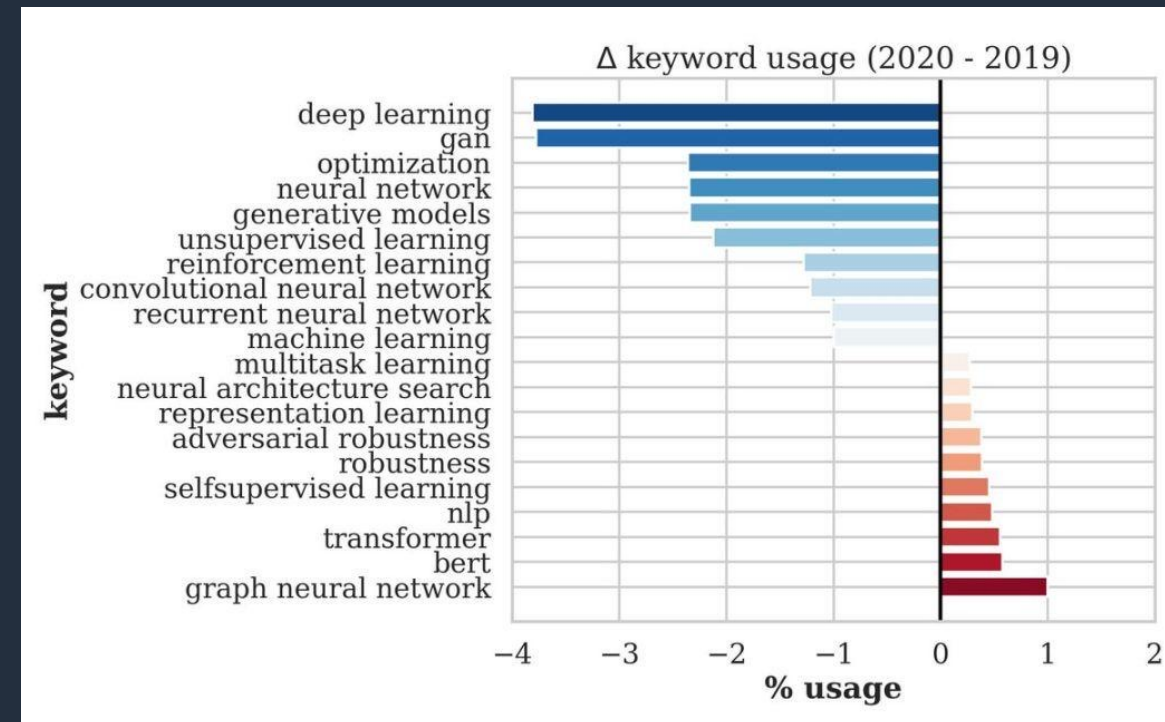
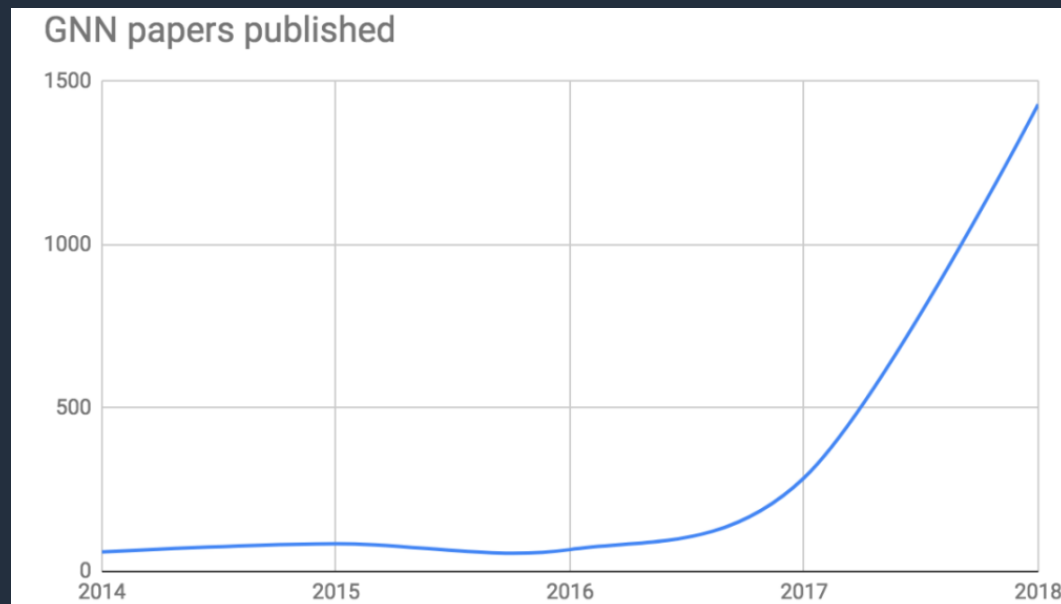
		Features																			
Nodes	1411	0	1	2	1	0	0	0	1	1	0	1	0	0	1	1	2	2	0	0	
	1410	0	1	1	1	0	1	0	0	1	0	1	0	1	0	1	1	1	1	1	
	338	0	0	0	0	1	0	1	0	0	1	0	0	0	1	0	0	0	0	0	
	339	1	0	0	0	2	0	1	0	0	2	0	1	0	1	0	1	0	0	0	
	1415	0	1	1	2	0	1	0	0	0	0	0	0	1	1	1	1	1	1	1	
	941	0	0	0	0	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0	
	1414	0	1	1	1	0	1	0	0	0	0	0	0	0	1	1	0	1	1	1	
	942	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	1413	0	1	1	1	0	1	1	0	0	0	0	0	0	1	1	0	1	1	1	
	1412	0	0	0	0	0	0	0	1	2	0	1	1	0	0	1	2	0	0	0	
	940	0	0	1	0	0	0	0	1	0	0	0	1	1	0	1	1	1	1	1	
	1419	0	0	1	0	0	1	0	1	1	0	1	1	1	1	0	1	1	1	1	
	945	0	1	4	3	0	0	0	0	2	0	1	0	0	2	1	0	3	1	1	
	332	0	0	0	0	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0	
	1418	0	0	1	0	0	0	0	1	0	0	0	1	2	0	1	0	1	0	1	
	946	0	1	1	0	0	1	0	1	0	0	0	1	4	0	1	1	1	2	0	
	333	0	0	0	0	1	0	1	0	0	1	0	0	0	0	0	0	0	0	0	
	1417	0	1	1	1	0	2	0	0	1	0	1	0	1	0	1	0	1	1	1	
	943	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	
	330	1	3	2	0	1	2	0	2	2	2	2	0	3	1	0	2	5	0	0	
1416	0	1	1	1	1	2	0	0	1	0	0	0	1	0	1	0	0	1	1		
944	0	1	4	2	0	0	0	2	0	1	0	0	2	0	0	3	1	0	0		
331	0	3	2	1	0	1	0	0	2	0	2	0	2	0	2	0	1	2	5		
949	0	0	0	0	2	0	0	1	0	1	0	1	0	0	0	0	0	0	0		
336	0	0	0	0	2	0	0	1	1	1	1	1	1	0	0	0	0	1	0		
337	1	1	1	0	0	1	2	0	1	1	1	0	1	1	1	1	1	1	1		
947	1	0	0	0	2	0	1	0	0	2	0	1	0	1	0	1	0	0	0		
334	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0		
948	0	0	0	0	1	0	1	1	0	1	1	0	1	1	1	0	1	1	0		
335	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0		
531	1	0	0	0	1	0	2	0	0	2	0	0	0	0	0	2	0	0	0		

# Graph Neural Networks (GNNs)

# Graph Neural Network (GNN)

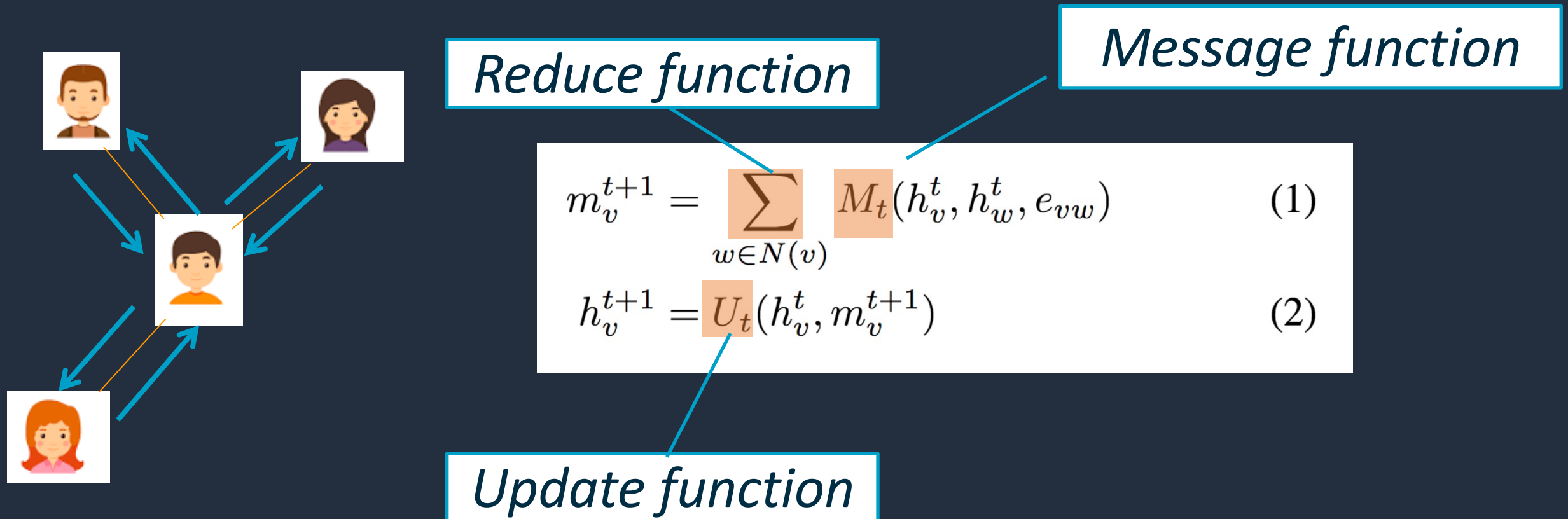
DL(NN) + Graph => GNN

A family of (deep) neural networks that learn node, edge, and graph features



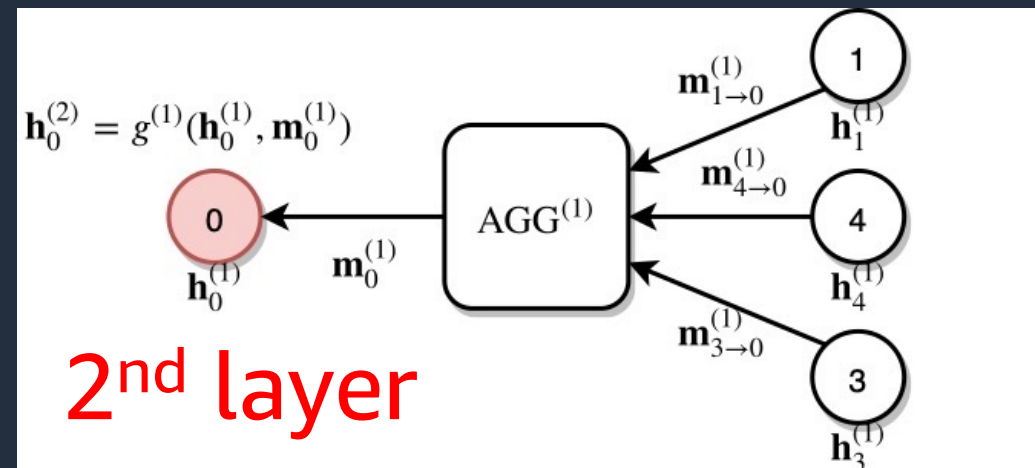
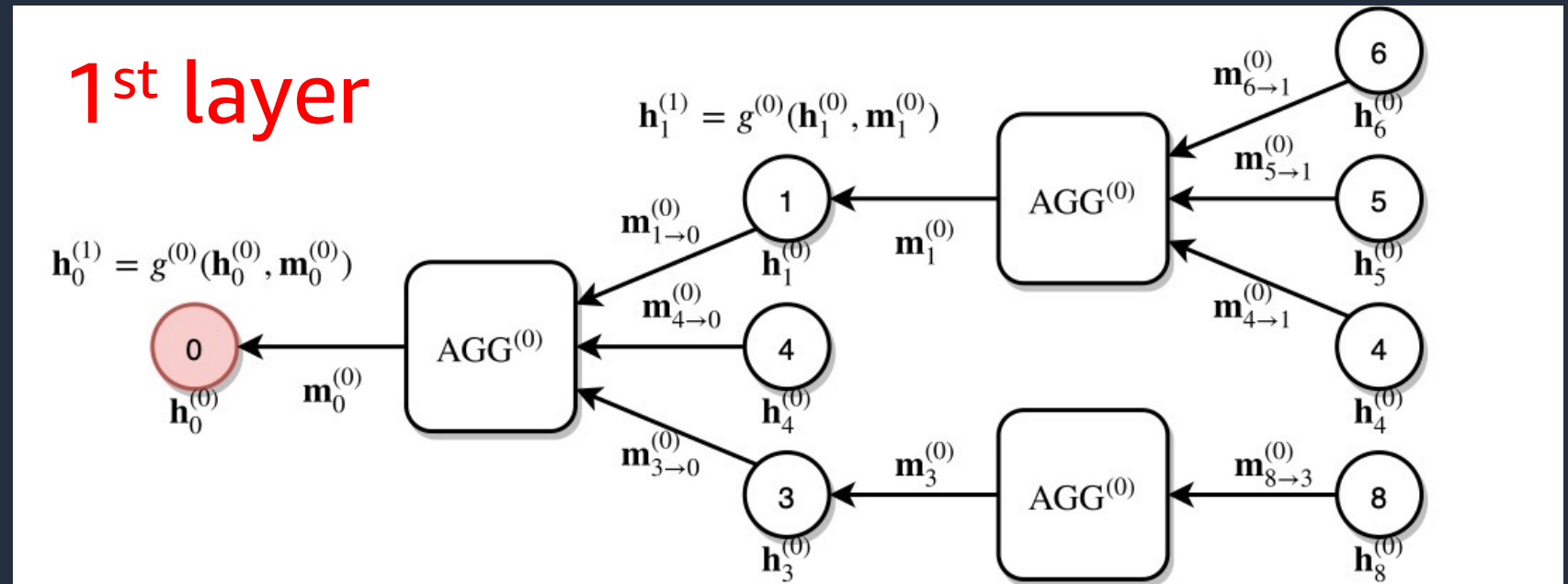
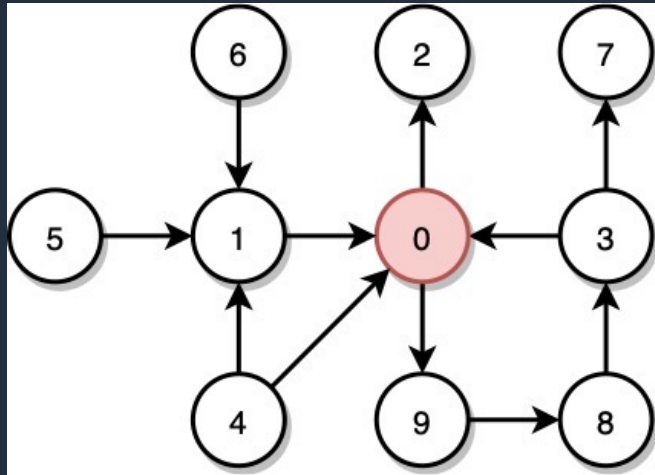
# How does GNN work

## Message-passing & Aggregation



# Stacked Multiple GNN layers

GNNs can *integrate* topologically distant information in a non-linear fashion.





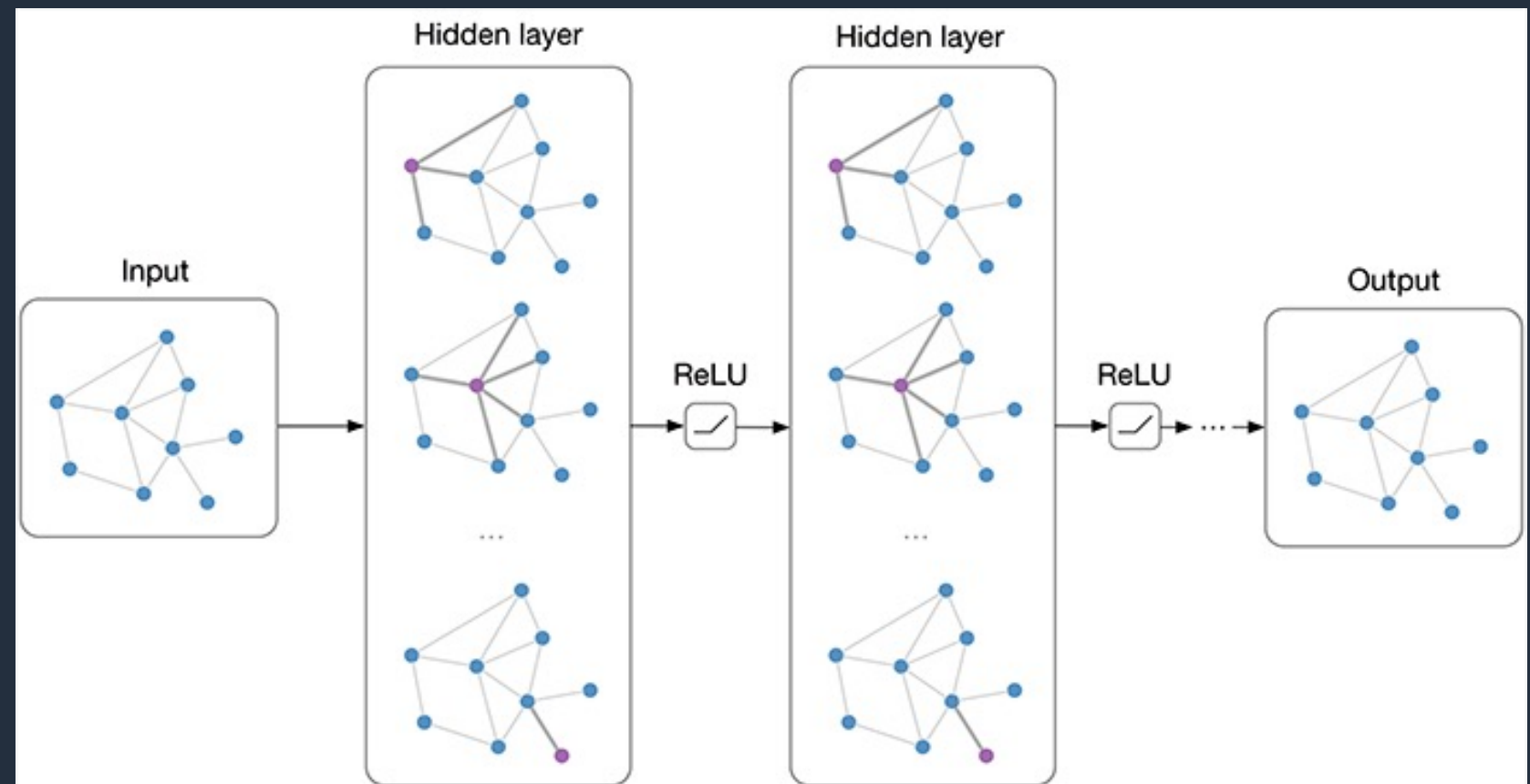
# Common GNN Models — GCN

GCN (Graph Convolutional Network) for homogeneous graphs

$$M_{vw}^{(l)} = \frac{h_w^{(l-1)}}{d_v + 1}$$

$$m_v^{(l)} = \sum_{w \in N(v) \cup \{v\}} M_{vw}^{(l)}$$

$$h_v^{(l)} = \phi(m_v^{(l)} W^{(l)})$$



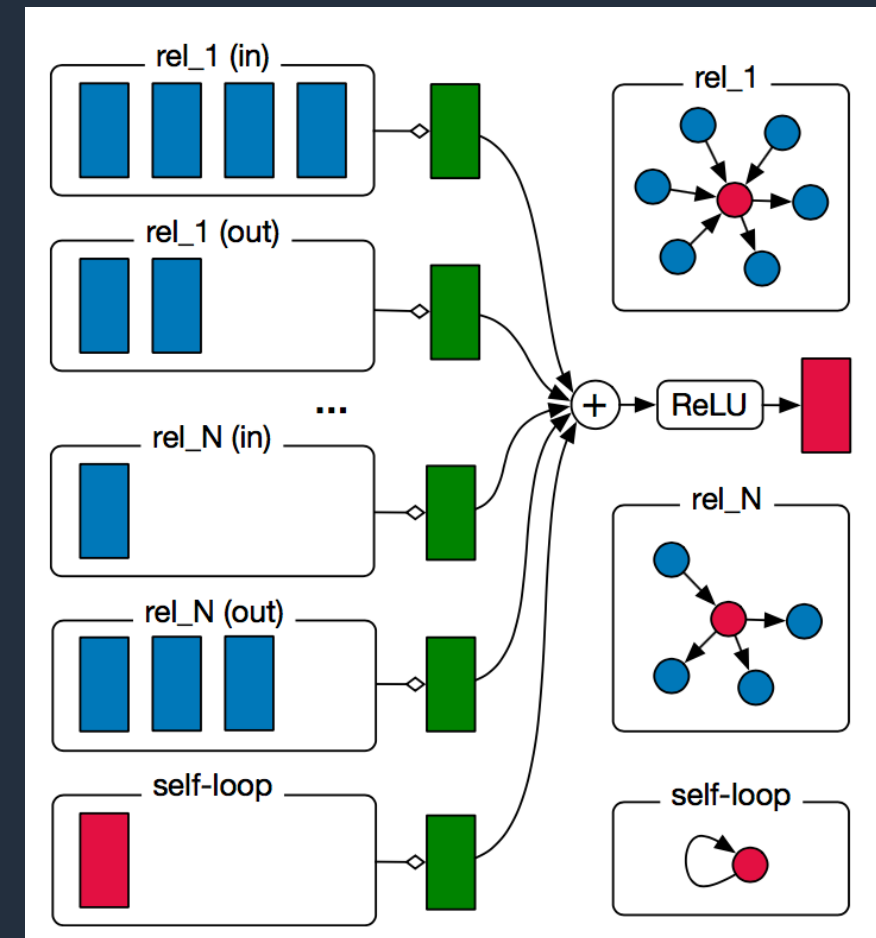
# Common GNN Models — RGCN

Relational graph convolution networks (RGCN) handles graphs whose nodes are connected with different relations.

$$M_{vw}^{(l)} = \frac{1}{c_{v,r}} W_r^{(l)} h_w^{(l-1)}, \text{ } r \text{ is the relation of } e_{vw}$$

$$m_v^{(l)} = \sum_{w \in N(v) \cup \{v\}} M_{vw}^{(l)}$$

$$h_v^{(l)} = \sigma(m_v^{(l)} W^{(l)})$$



# Common GNN Models — GAT

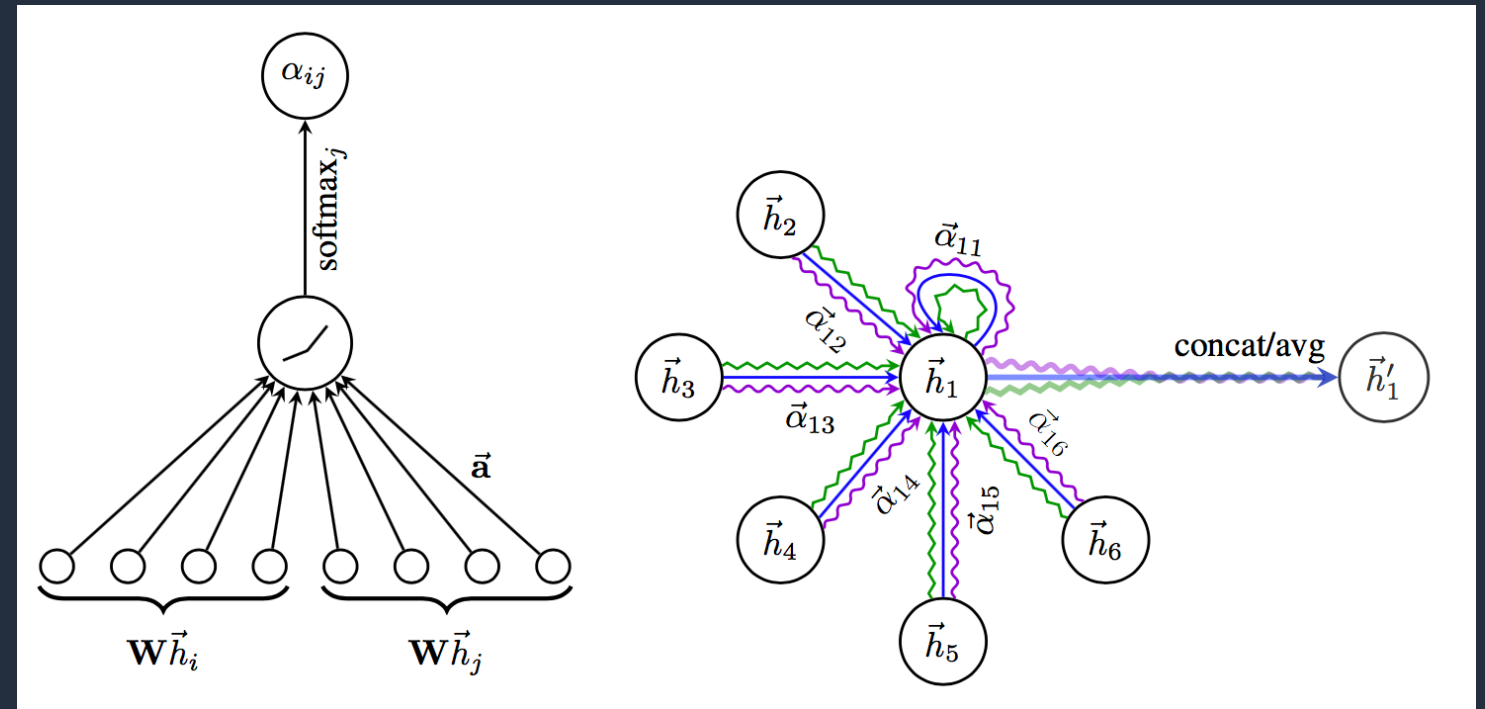
Graph Attention Network(GAT) provides weighted sum over the neighborhood—  
Enables to selectively integrate information.

$$M_{vw}^{(l)} = \alpha_{vw} h_w^{(l-1)}$$

$$m_v^{(l)} = \sum_{w \in N(v) \cup \{v\}} M_{vw}^{(l)}$$

$$h_v^{(l)} = \phi(m_v^{(l)} W^{(l)})$$

$$\alpha_{vw} = \frac{\exp(\text{LeakyReLU}(\vec{a}^T [W \vec{h}_v || W \vec{h}_w]))}{\sum_{k \in N_v} \exp(\text{LeakyReLU}(\vec{a}^T [W \vec{h}_v || W \vec{h}_k]))}$$



# GNN References

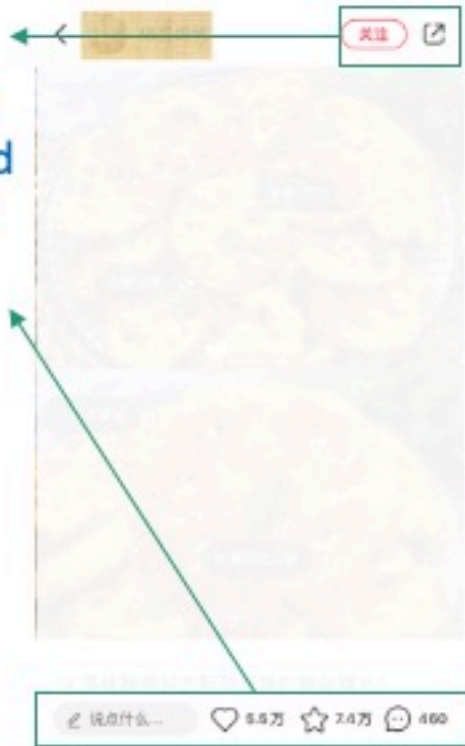
- GNN Libraries
  - DGL (Deep Graph Library): <https://www.dgl.ai>
  - PyG (Pytorch Geometric): <https://pyg.org/>
  - TF\_GNN (TensorFlow GNN): <https://github.com/tensorflow/gnn>
- Online Books
  - Deep Learning on GraphS ([https://yaoma24.github.io/dlg\\_book/](https://yaoma24.github.io/dlg_book/))
  - Graph Neural Networks (<https://graph-neural-networks.github.io/>)
- Online Courses
  - Stanford CS224W: ML with Graphs(<https://web.stanford.edu/class/cs224w/>)

# Real Cases of Using GNN

# GNN Helps to Detect 'Suspicious' Interaction

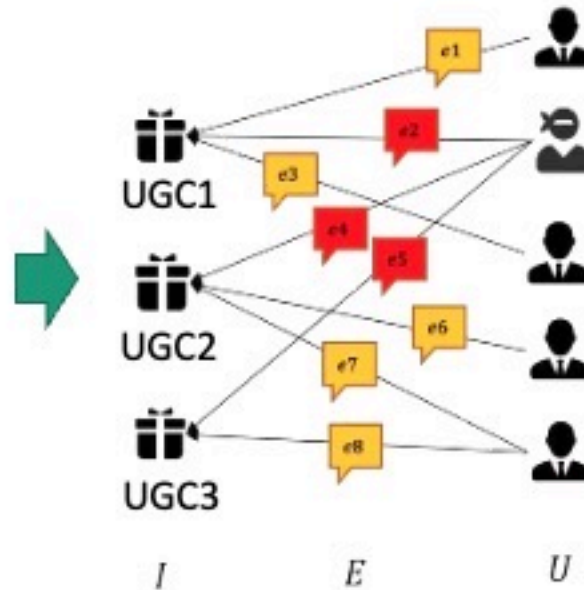
## Problem:

Interactions between users and UGCs could be 'suspicious' with hidden agendas.

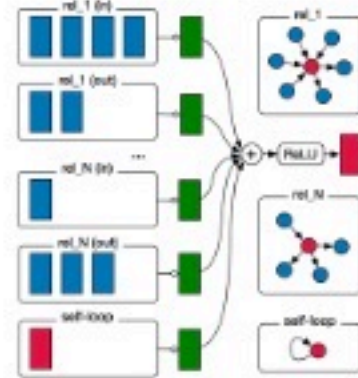


## Bipartite:

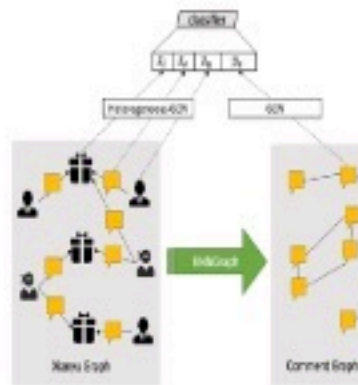
User  $\leftrightarrow$  UGCs



## Improvements:



Focus on User:  
As a node classification task, RGNN model out-performs baseline models in all settings



Focus on Behavior:  
SOTA model out-perform baseline models in some settings, and help find **new patterns** of 'suspicious' that previous rule-based methods can not touch.

# GNN Helps to Detect 'Bot' Accounts

## Problem:

Bots are increasingly impacting the e-comm platform's customer in:

- promo code abuse.
- inventory encumbrance.
- return logistics costs.
- fraud response operations cost.

## Challenges:

- Account features differ in different life stages.
- New/Inactive accounts having nearly no features.
- Lack of solid labels, particularly newly registered.
- Huge amounts.

## Tables -> Bipartite:

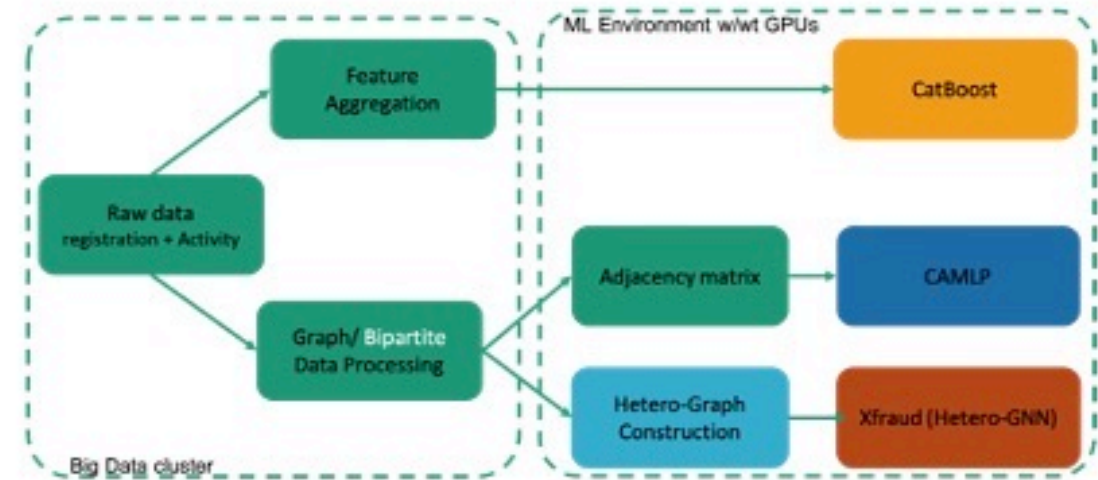
UID	isBot	Time	IP	App	version	city	...
1	0	86400	59.102.48.1	app1	6.9.2	BJ	...
2	1	86401	176.201.138.67	app2	4.1	BJ	...
...	...	...	...	...	...	...	...
N	...	...	...	...	...	...	...

UID	Platform	version	Phone	...
1	iOS	12.1	010-694899309	...
2	iOS	12.10	025-084850598	...
3	Android	9.5	085-1340586763	...
...	...	...	...	...
N	Android	9.5	010-445648935	...

UID	IP	UID	App_version
1	59.102.48.*	1	App1 6.9.2
2	176.201.138.*	2	App2 4.1
...	...	...	...
N	95.218.200.*	N	Web 5.5.0



## Improvements:

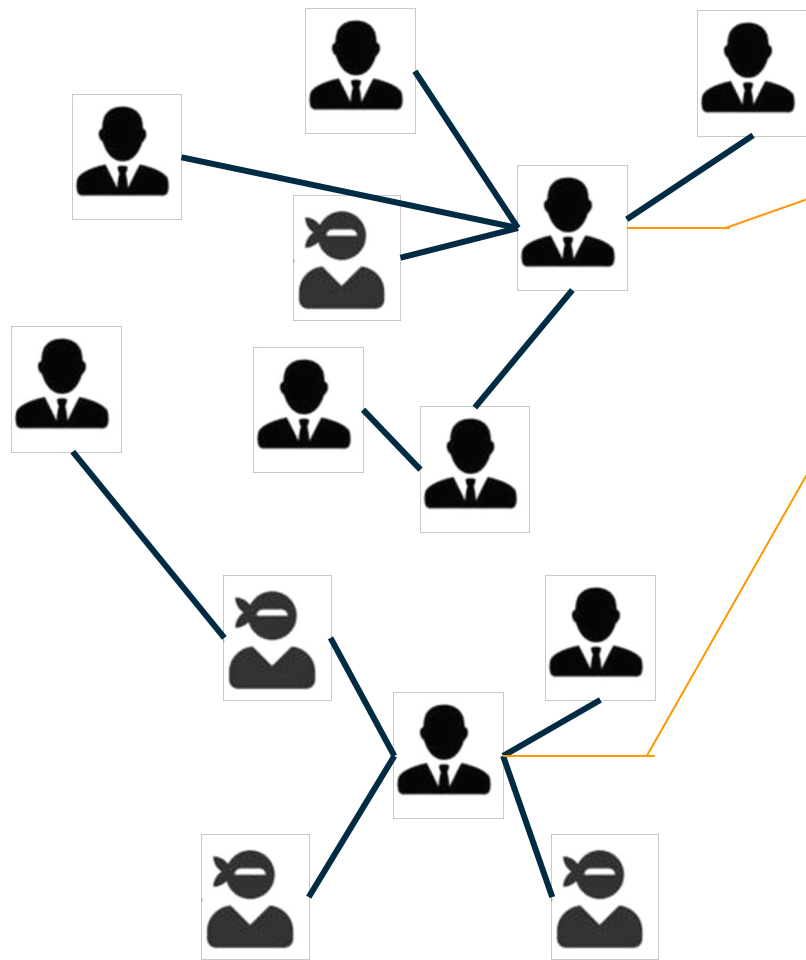


Multiple models all have 11x performance boost than customer's existing models!  
GNN and tree-based models work together, helping to handle accounts in their whole-life.

# Case 3: A FinTech – Predict Loan Overdue



# Problem: Personal loan overdue prediction

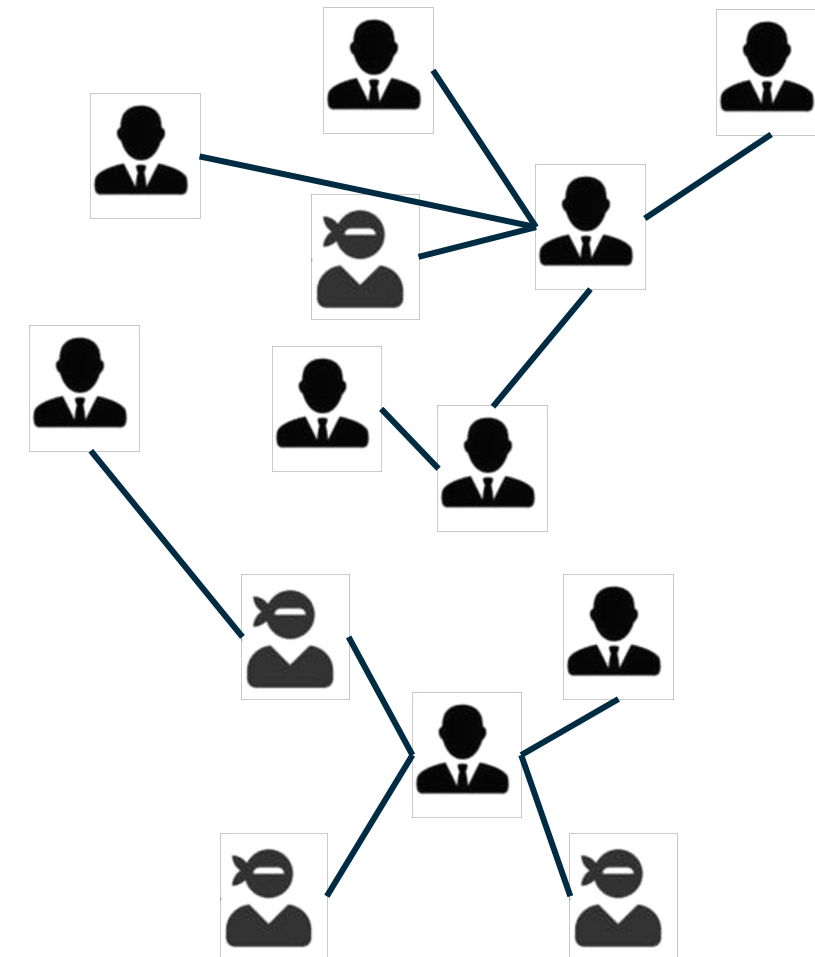


Will the personal loan application be overdue?

- Has an xGboost model as baseline, using person's info, loan history, and social relation as input features
- Hope to fully leverage the social relation about customers to boost prediction performance

# The Data - Homogenous graph

Item	Performance Verification
Data range*	7 years
No. Nodes – phone*	26M
No. Edges – contact*	208M
Node features	353 features
Negative nodes	2.3M
Positive nodes	16K

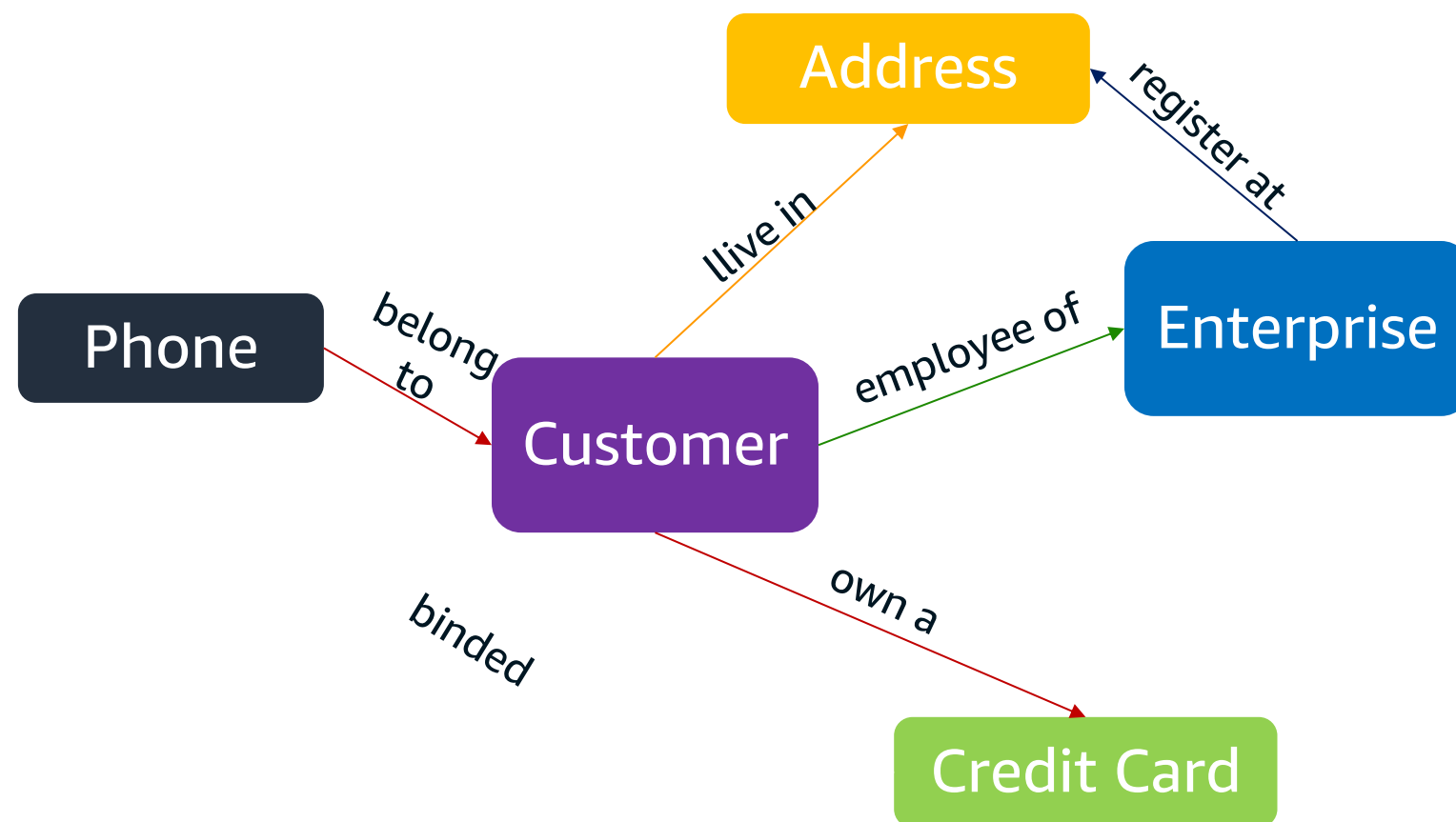


# Model Selection — Node Classification

- Most of GNN Models do node classification
  - GCN (Graph Convolutional Network)
  - GraphSage (Simplified GCN)
  - Graph Attention Network (GAT)
- Pre-extract higher-order local structure feature
  - Reflex, GDV (Graphlet Degree Vector), etc.
- GraphSage model out-performance GCN and GAT, achieving +4% AUC than the xGboost baseline.

# Case 4: A Bank – Loan Application Approval

# Problem: Loan application approval



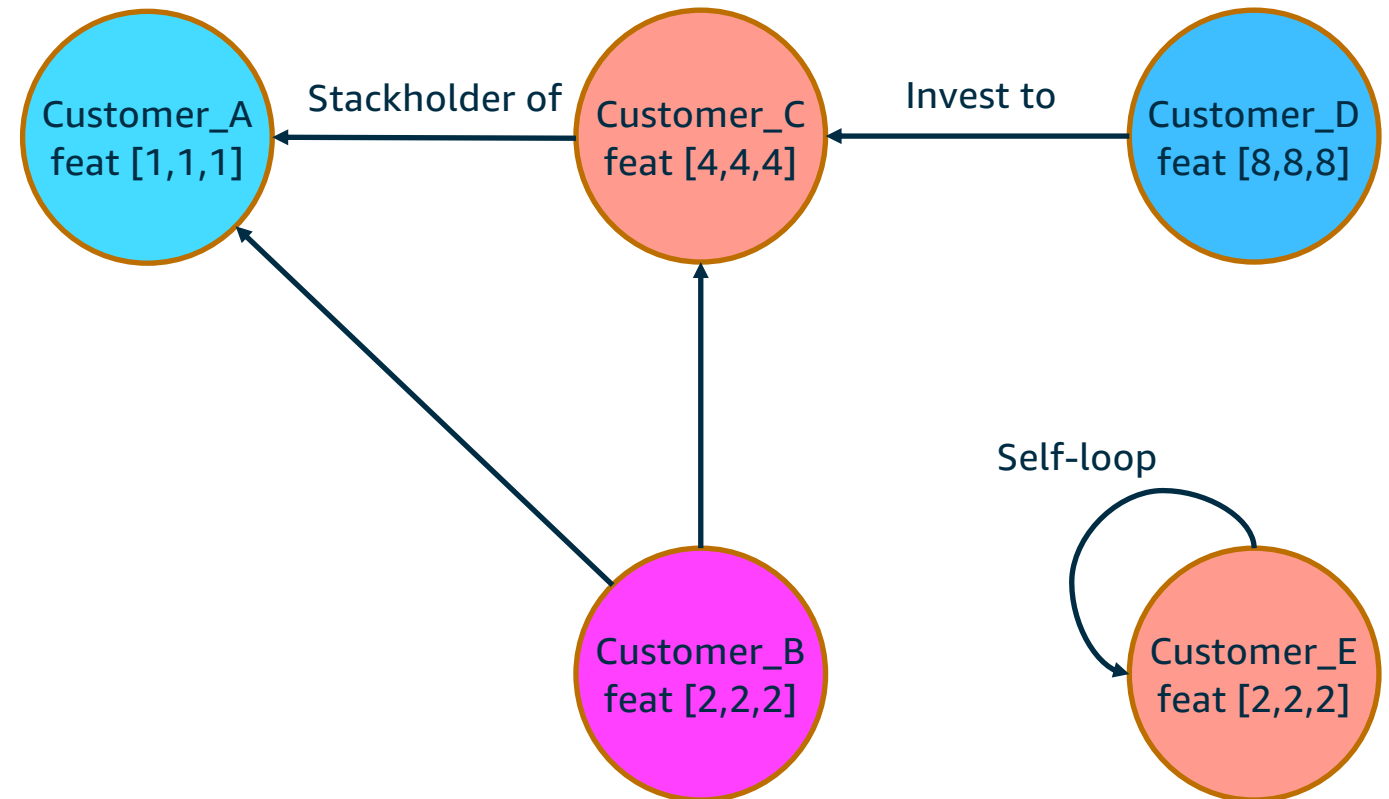
Knowledge Graph in  
Banking (demo)

The bank has built a knowledge graph about many entities and relations, which could help to predict risk scores for various business scenarios

- Have man-made rules for loan approval decision
- Hope to leverage GNN models to complete the approval in real time

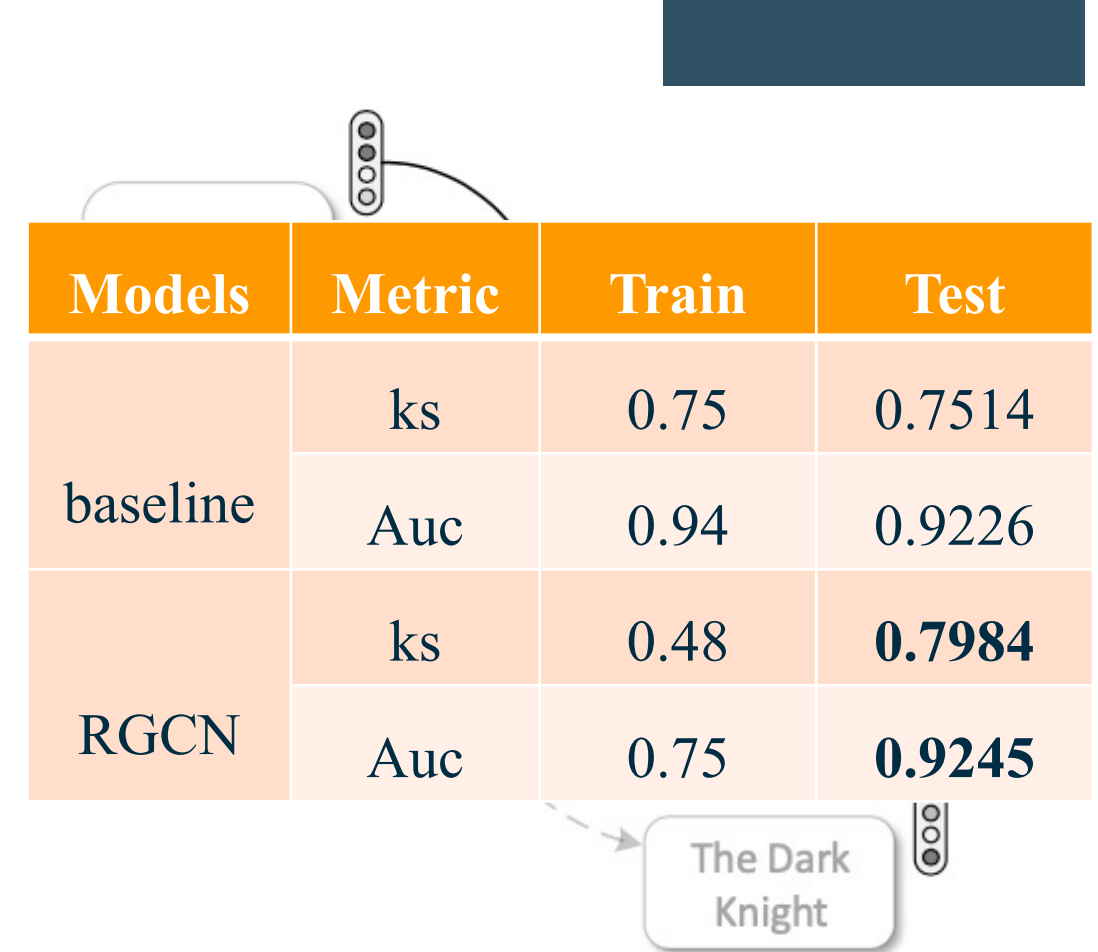
# The Data – Knowledge Graph (KG)

- The Graph
  - 2 types of entities
  - ~20 types of relationships
  - contents coming from various sources, e.g. banks, governments, and etc.
- Labels
  - Credit card users' payment history
  - Credit scores from government agencies.



# Model Selection — KG-oriented

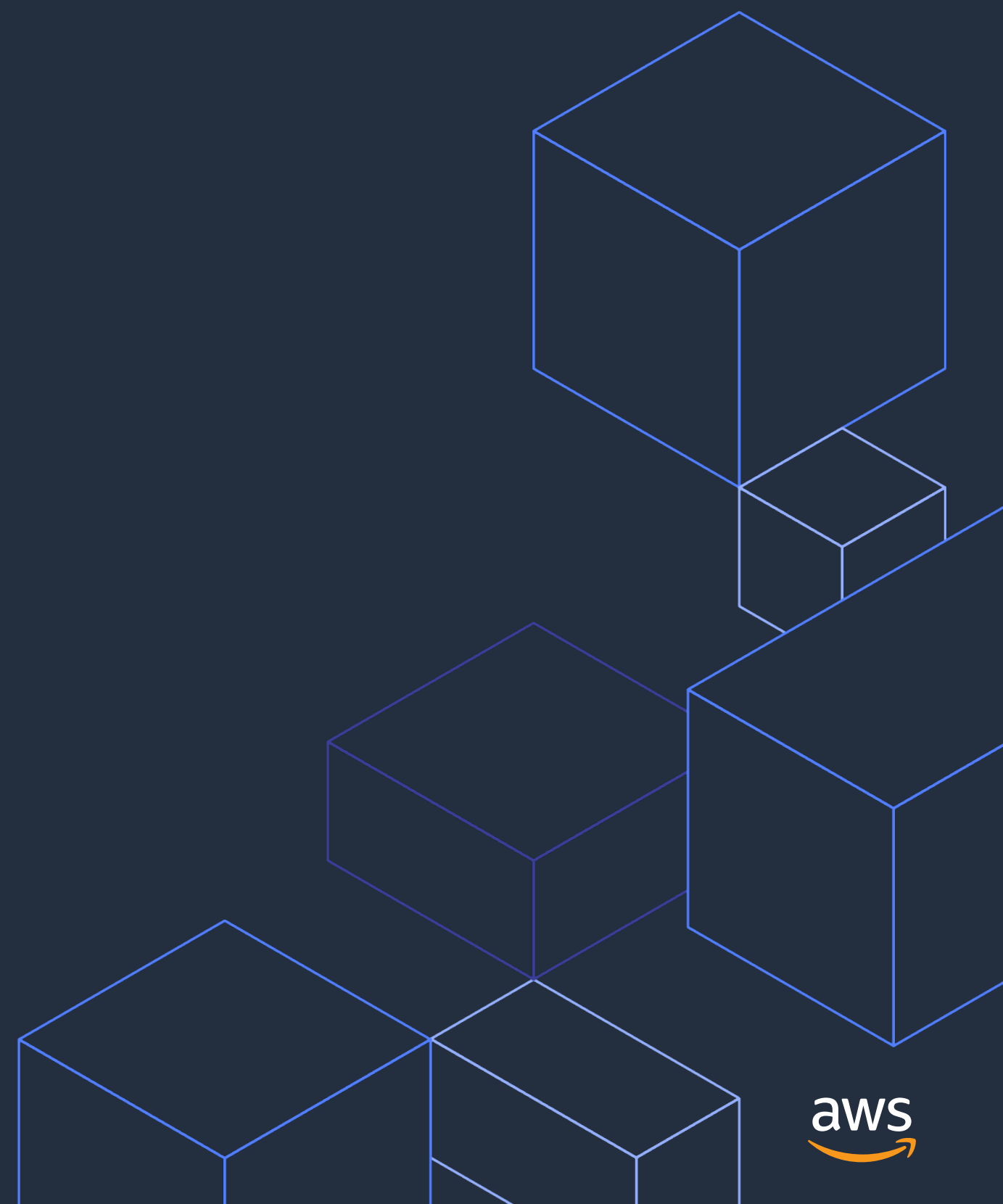
- RGCN
  - KG is a specious case of heterogenous graph, where RGCN can work on too
  - Prone to overfitting if the no. of relations is large.
- ICLR 2020 — CompGCN
  - Specially designed for KG
  - Has less parameters, reported SOTA performance



Models	Metric	Train	Test
baseline	ks	0.75	0.7514
	Auc	0.94	0.9226
RGCN	ks	0.48	<b>0.7984</b>
	Auc	0.75	<b>0.9245</b>

The Dark Knight

# GraphStorm





# Challenges of adopting GNN



**Steep Learning  
Curve**



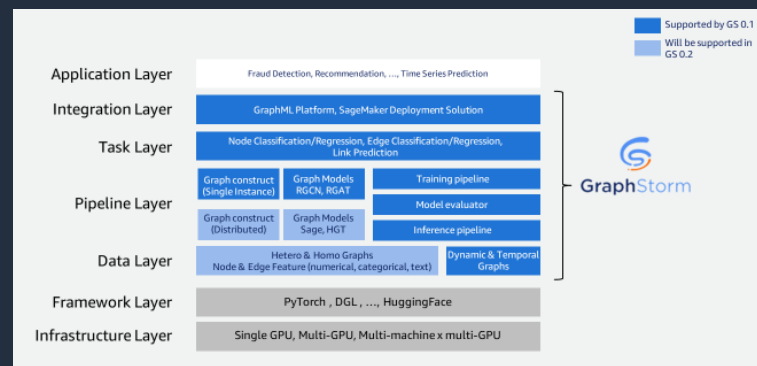
**Complex graph  
data processing**



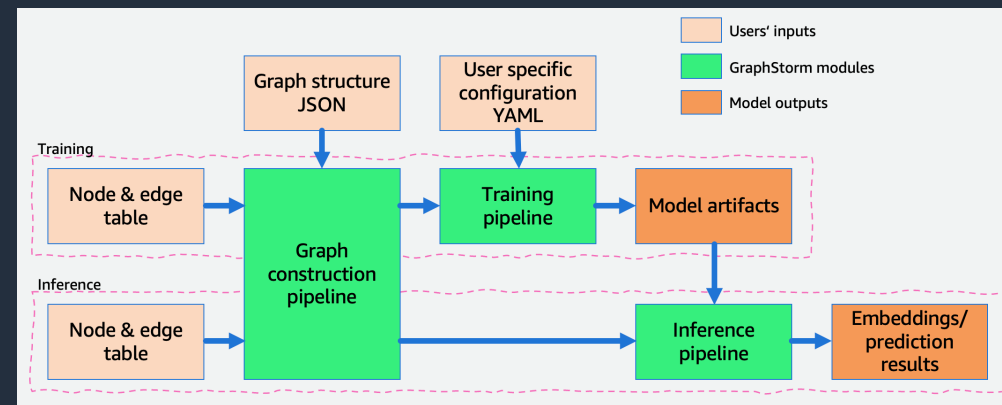
**Hard to Scale to  
extreme large graphs**

# Fast-track graph ML with GraphStorm: A new way to solve problems on enterprise-scale graphs

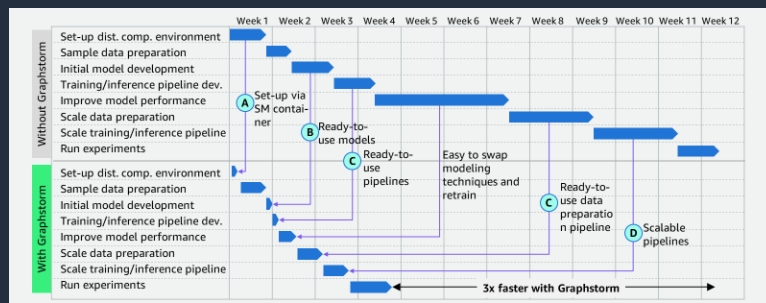
## Comprehensive toolbox



## Easy-to-use interface



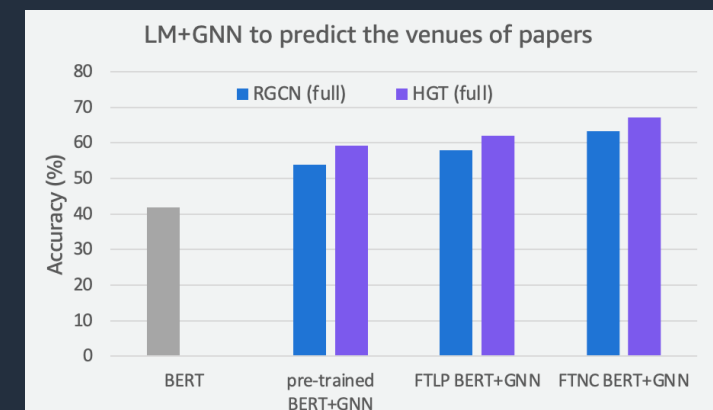
## Speed up model development & deployment



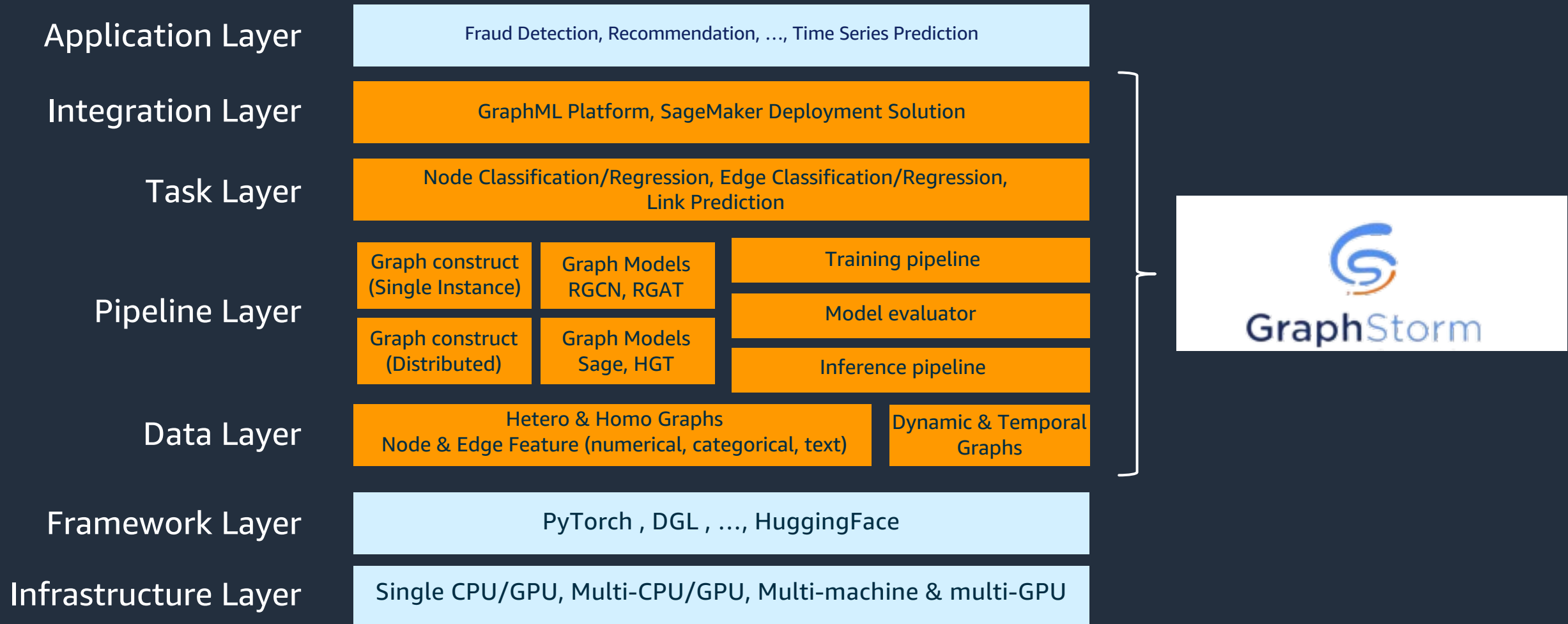
## Scale to billion-node graphs



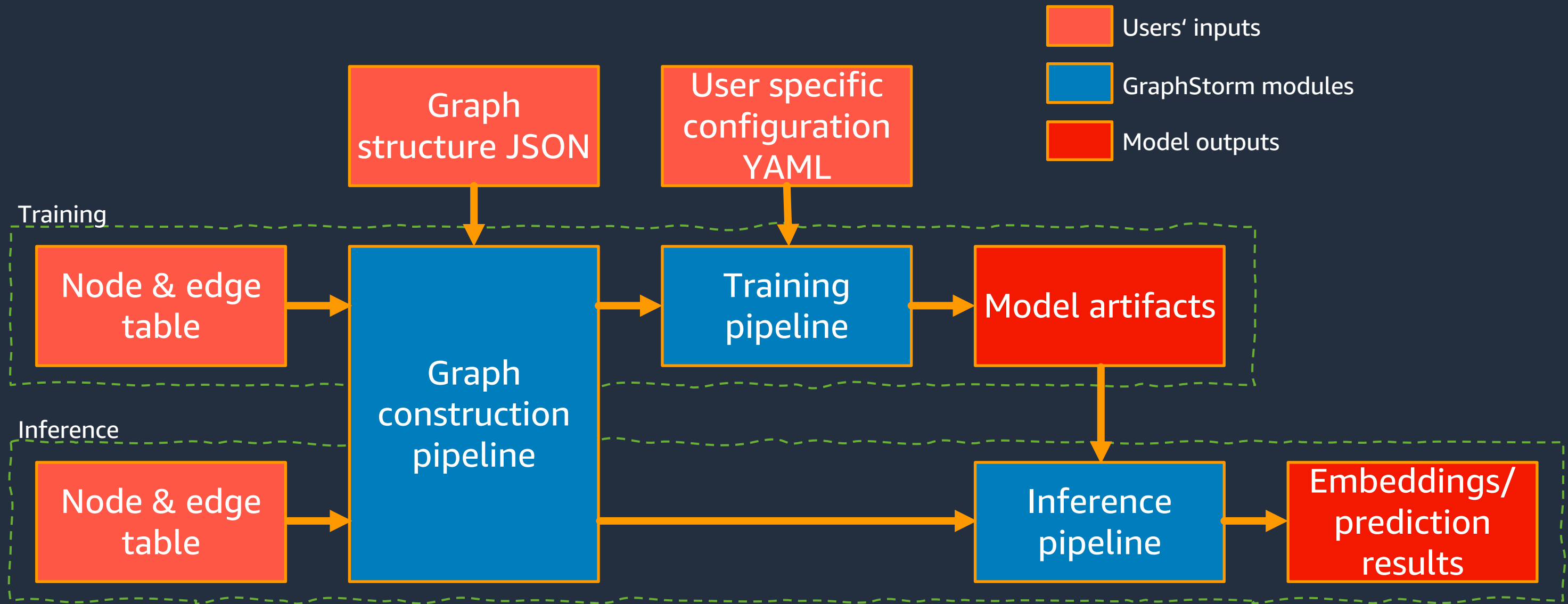
## Superb model performance



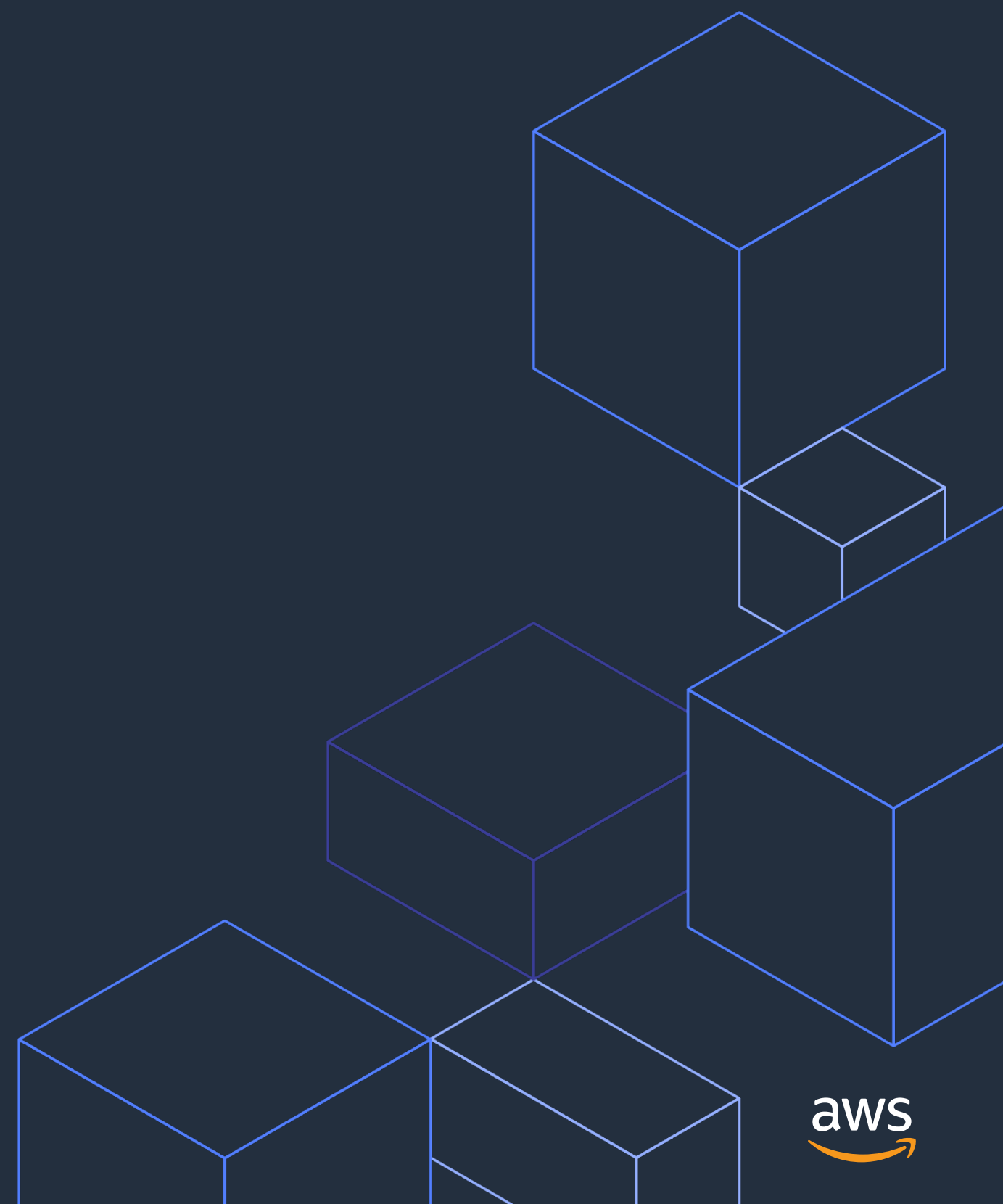
# Comprehensive: many functionalities out-of-the-box



# Easy-to-use: ready-to-use pipelines



# 15min Break



# GraphStorm hands-on



# Hands-on Environment

Each one will have a temporary account to create an AWS EC2 instance via:

<https://catalog.us-east-1.prod.workshops.aws/join?access-code=e94d-04a4a4-d3>

# Q&A