# GNN Advanced
## -- GNN in NLP
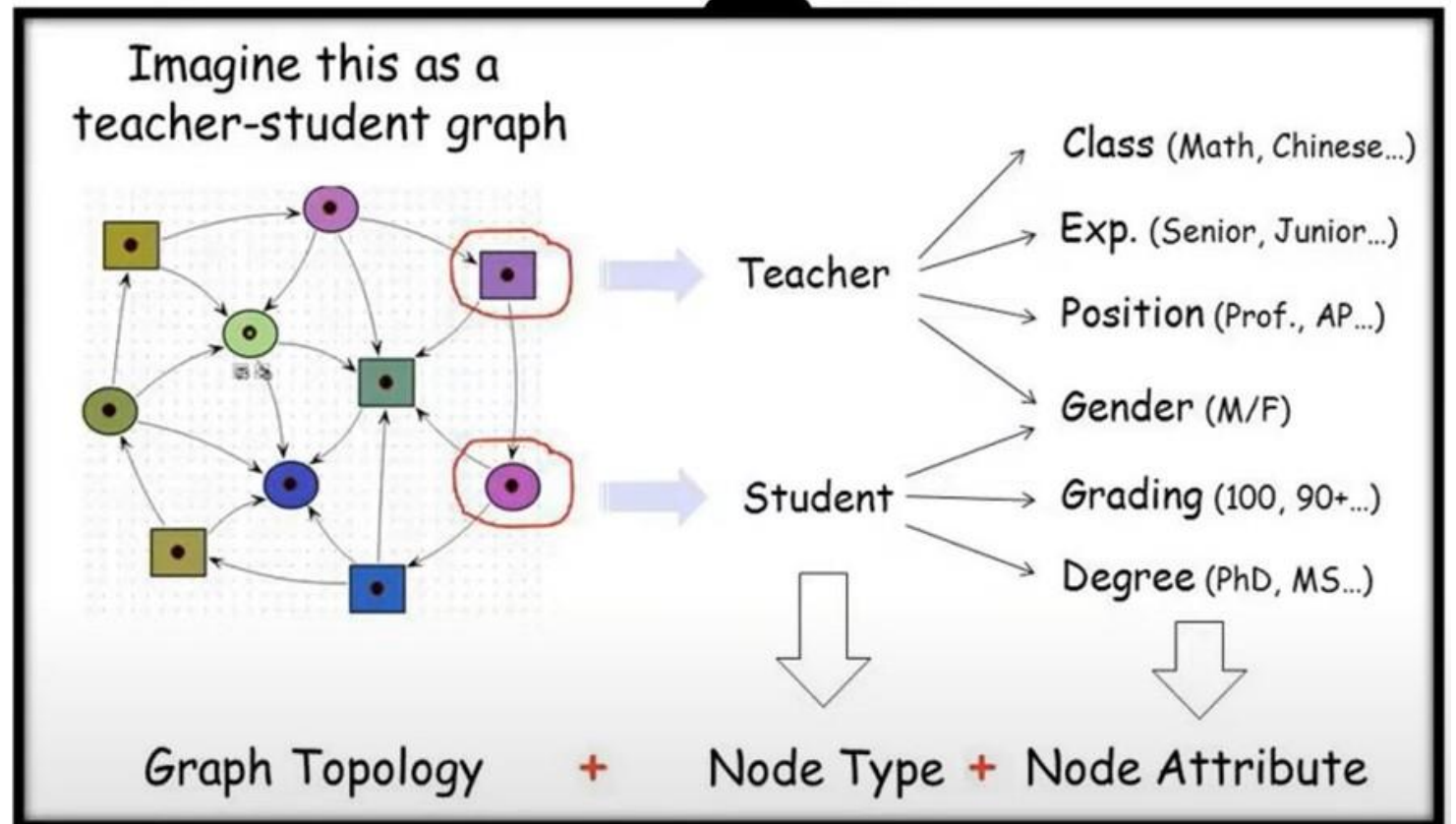
Zhaoxuan Tan

Xi'an Jiaotong University

March 5, 2022

# Table of Contents

- Graph Construction
  - Static Graph Construction
  - Dynamic Graph Construction
- Heterogeneous Graph Representation Learning
  - R-GCN
  - HAN
  - HGT
- GNN for Knowledge Graph Embedding
  - R-GCN / SACN / KBGAT / CompGCN
- Application
  - Syntactic GCN for SRL
  - RE-SIDE
  - QA-GNN

# Why Graph?

- Graphs are a general language for describing and modeling complex system
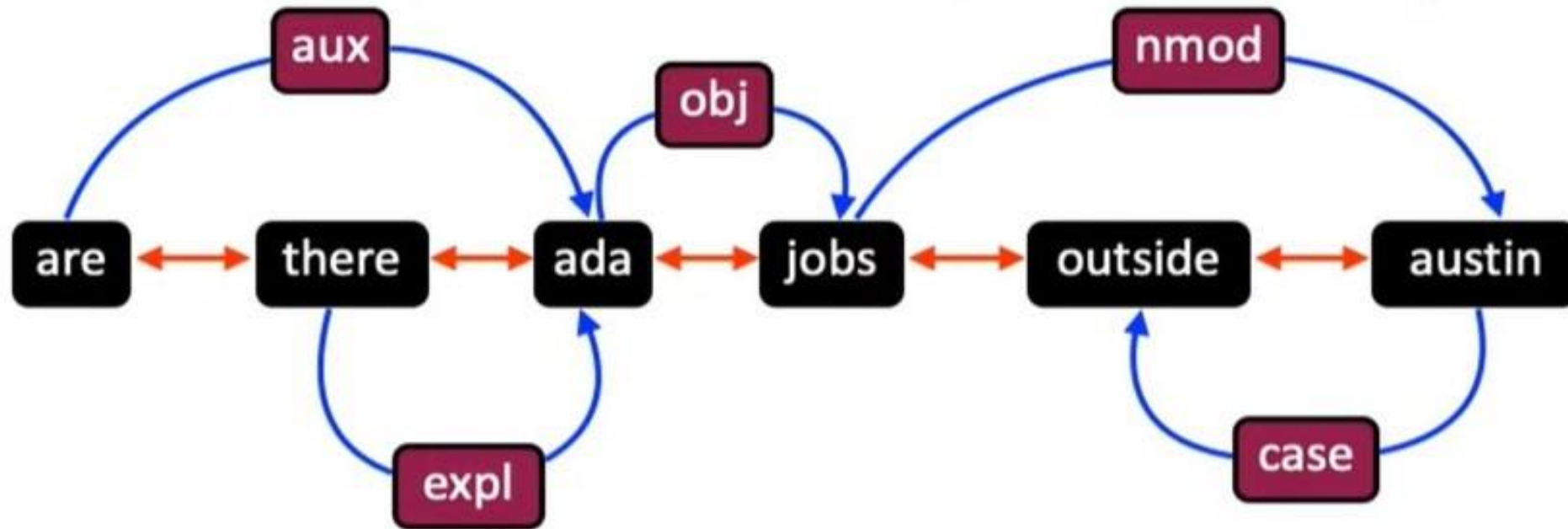
# Graph Construction: Static Method

- Problem Setting:
  - Input: raw text (sentence, paragraph, document etc.)
  - Output: graph

- Conducting during preprocessing by augmenting text with domain knowledge
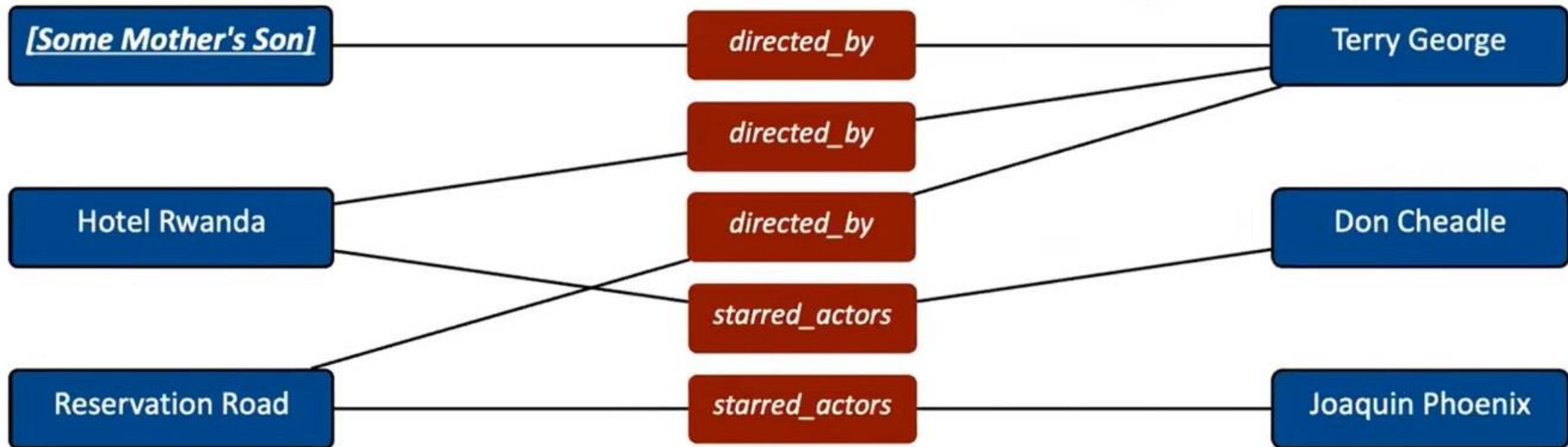
# Static Method: Dependency Graph

- Text: are there ada jobs outside Austin?



- Add sequential edge
  - reserve sequential information in raw text
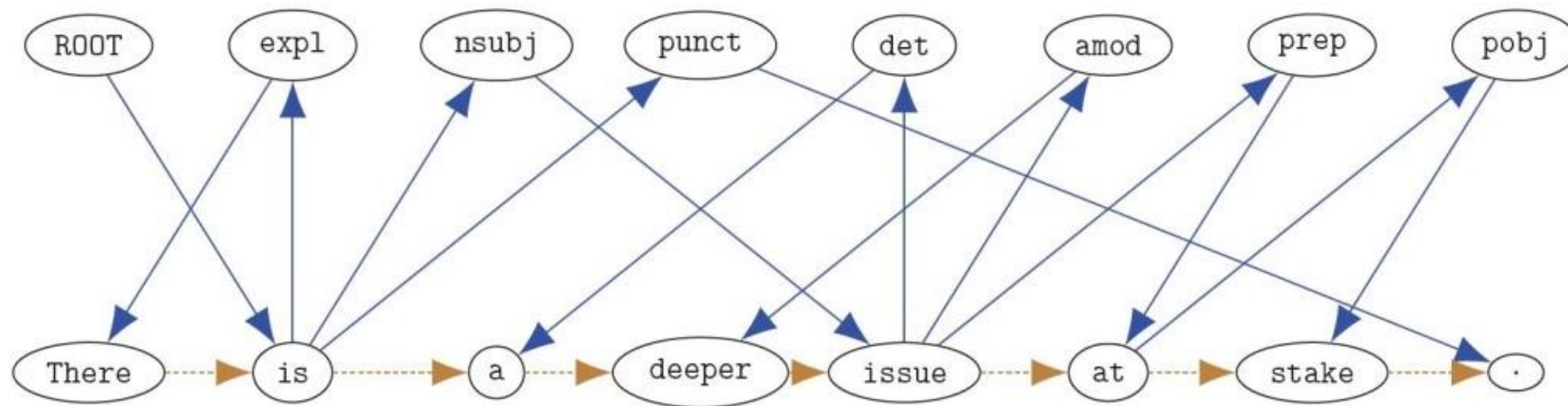  - Connection multiple dependency graphs in a paragraph

# Static Method: Knowledge Graph

- Question: who acted in the movie directed by the director of [*Some Mother's Son*]
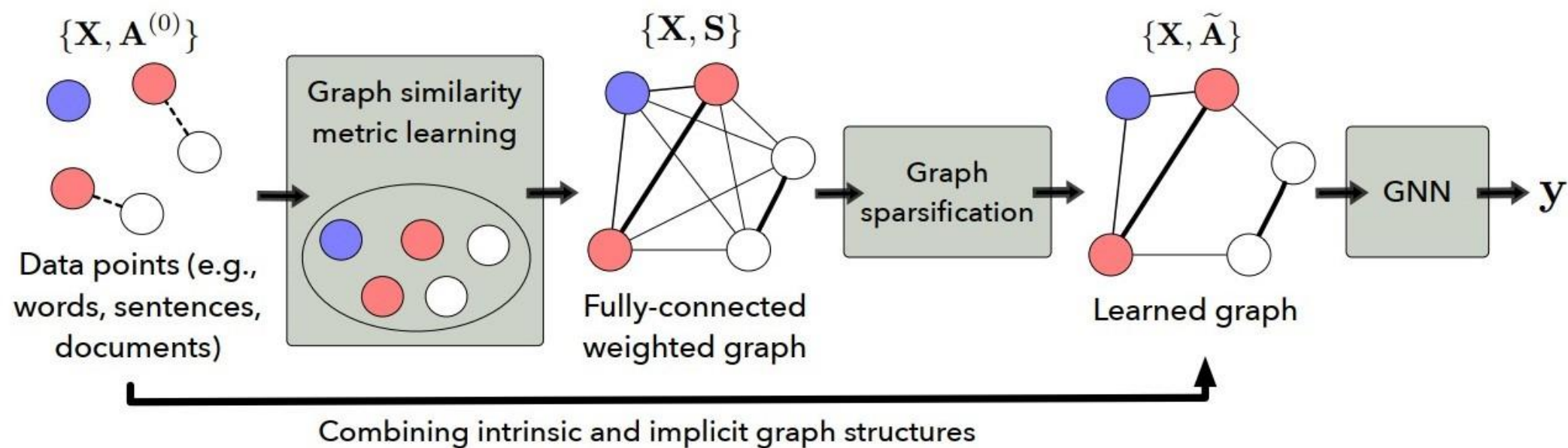- Answer: Don Cheadle, Joaquin Phoenix

# Static Method: Levi Graph

- An edge for every (node, edge)
- Edge label have hidden embedding

# Graph Construction: Dynamic Method



$\{\mathbf{X}, \mathbf{A}^{(0)}\}$

Data points (e.g., words, sentences, documents)

Graph similarity metric learning

$\{\mathbf{X}, \mathbf{S}\}$

Fully-connected weighted graph

Graph sparsification

$\{\mathbf{X}, \tilde{\mathbf{A}}\}$

Learned graph

GNN

$\mathbf{y}$

Combining intrinsic and implicit graph structures

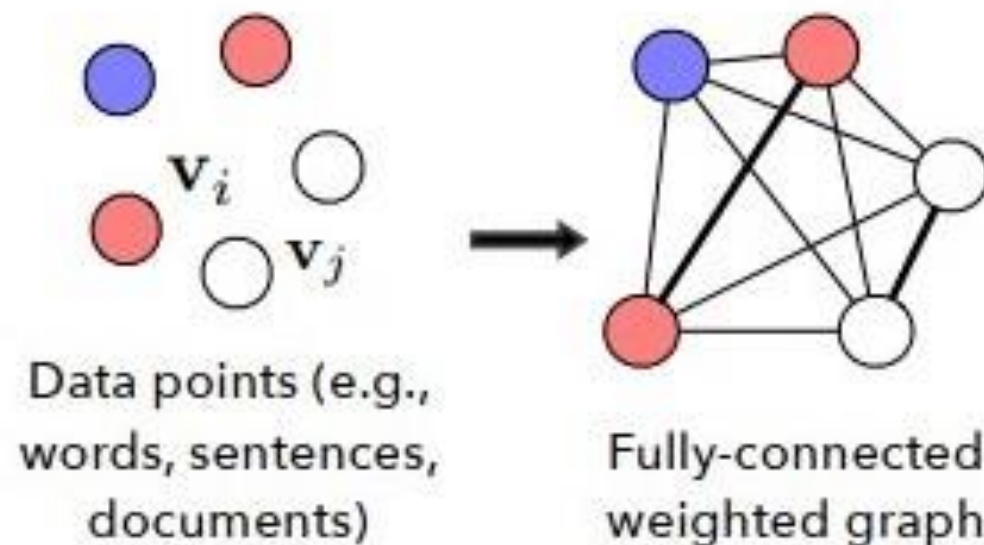# Dynamic Method: Attention-based node similarity

$$S_{i,j} = (\mathbf{v}_i) \odot (\mathbf{u})^T \mathbf{v}_j$$

Node feature vector

Non-negative learnable weight vector

$$S_{i,j} = \text{ReLU}(\mathbf{W}\mathbf{v}_i)^T \text{ReLU}(\mathbf{W}\mathbf{v}_j)$$
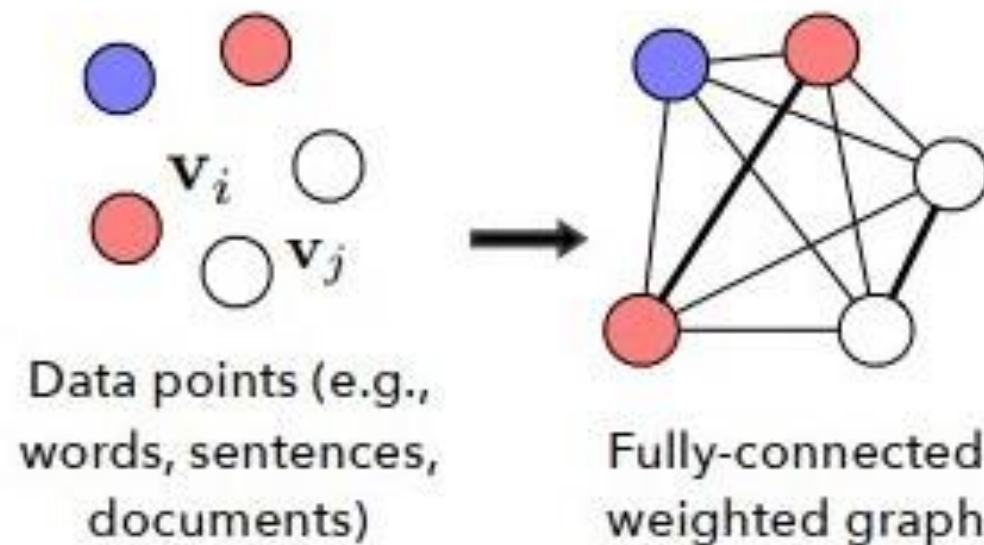
Learnable weight matrix

Data points (e.g., words, sentences, documents)

Fully-connected weighted graph

# Dynamic Method: Cosine-based node similarity

$$S^p_{i,j} = \cos(\mathbf{w}_p \odot \mathbf{v}_i, \mathbf{w}_p \odot \mathbf{v}_j)$$

Learnable weight vector

$$S_{i,j} = \frac{1}{m} \sum_{p=1}^{m} S^p_{ij}$$ Multi-head similarity scores

$\mathbf{v}_i$

$\mathbf{v}_j$

Data points (e.g., words, sentences, documents)

Fully-connected weighted graph

# Dynamic Method: Graph Sparsification
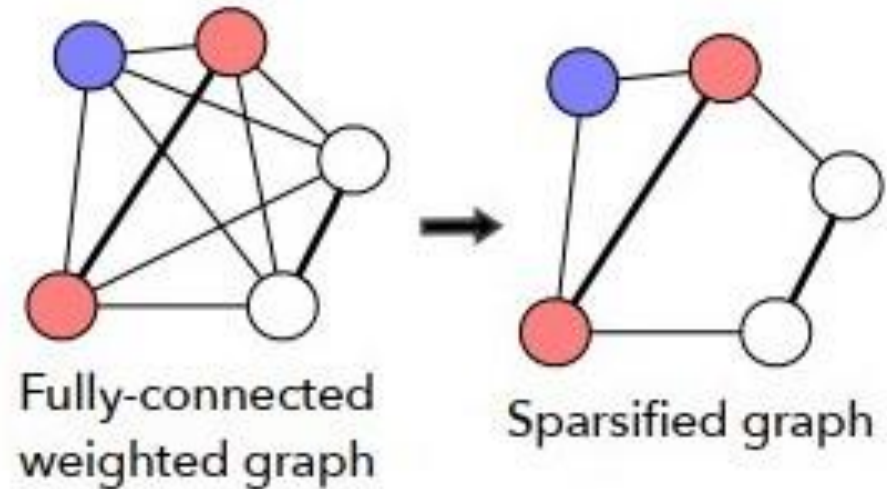
- KNN-style Sparsification

$$\mathbf{A}_{i,:} = \mathrm{topk}(\mathbf{S}_{i,:})$$

- Epsilon-neighborhood Sparsification

$$A_{i,j} = \begin{cases} S_{i,j} & S_{i,j} > \varepsilon \\ 0 & \text{otherwise} \end{cases}$$
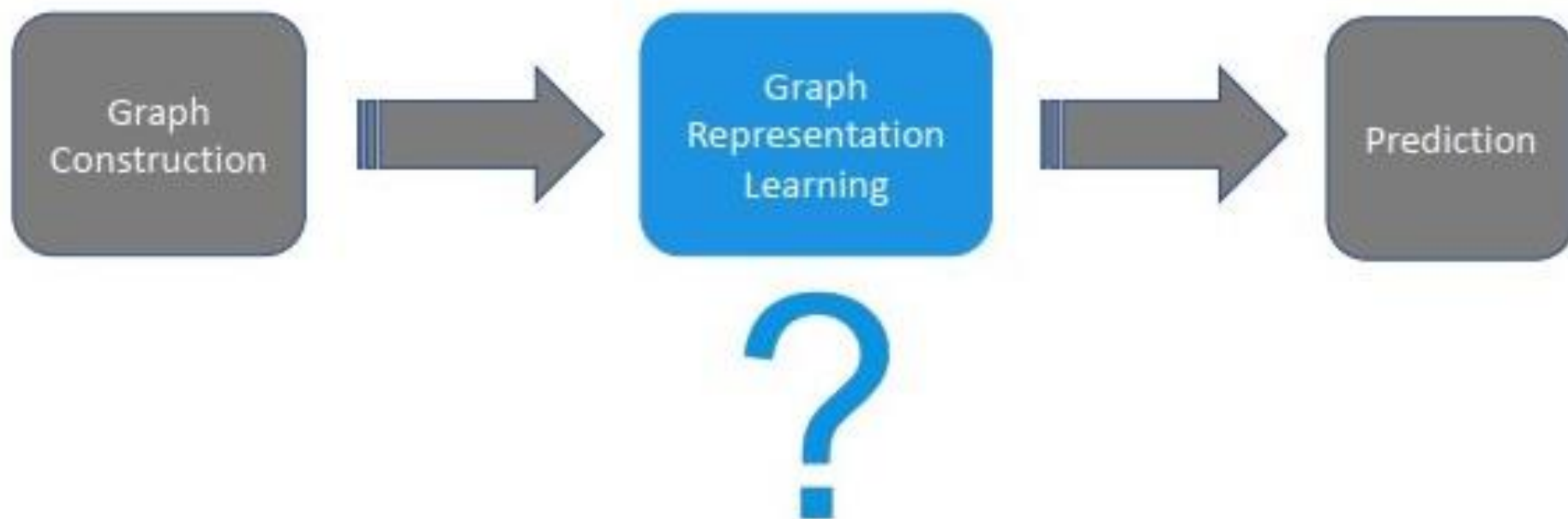
- Graph Regularization

$$\frac{1}{n^2}\|A\|_F^2$$



Fully-connected weighted graph

Sparsified graph

# Static vs. Dynamic Graph Construction

| Static graph construction | Dynamic graph construction |
|---|---|
| Pros | Pros |
| prior knowledge | no domain expertise |
| | joint graph structure & representation learning |
| Cons | Cons |
| extensive domain expertise | scalability |
| • error-prone (e.g., noisy, incomplete)<br>• sub-optimal | explainability |
| • disjoint graph structure & representation learning<br>• error accumulation | |

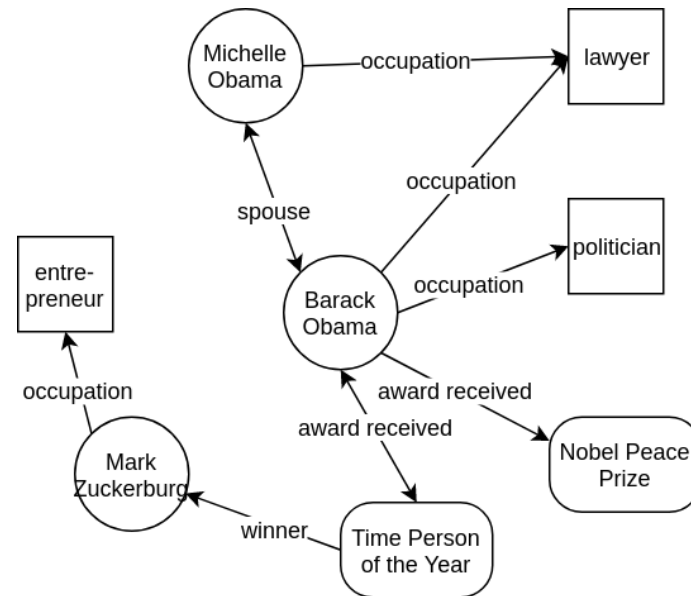# After constructing graph? -> Representation!
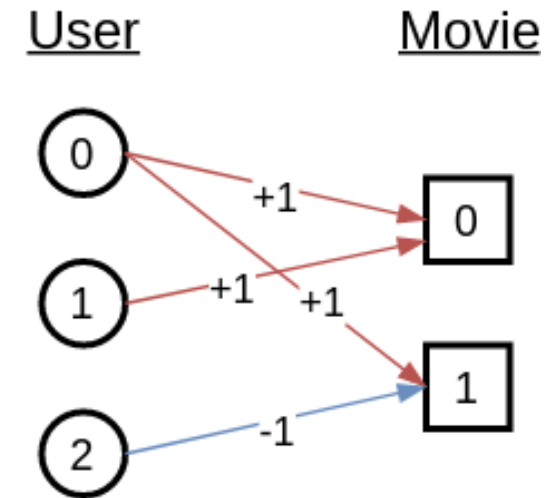
# What is Heterogeneous GNN?

- Citation Graph
- Knowledge Graph
- Recommender System
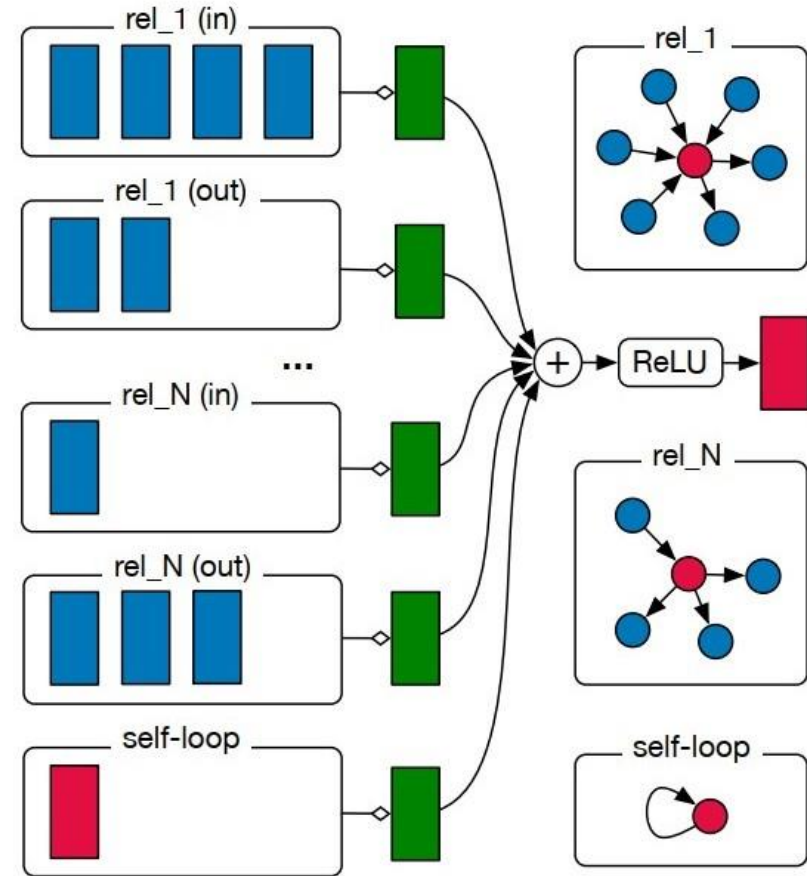
# R-GCN

- Message Passing Scheme

$$h_i^{(l+1)} = \sigma \left( \sum_{r \in \mathcal{R}} \sum_{j \in \mathcal{N}_i^r} \frac{1}{c_{i,r}} W_r^{(l)} h_j^{(l)} + W_0^{(l)} h_i^{(l)} \right)$$

- Regularization: avoid over-parameterization
  - Basis decomposition

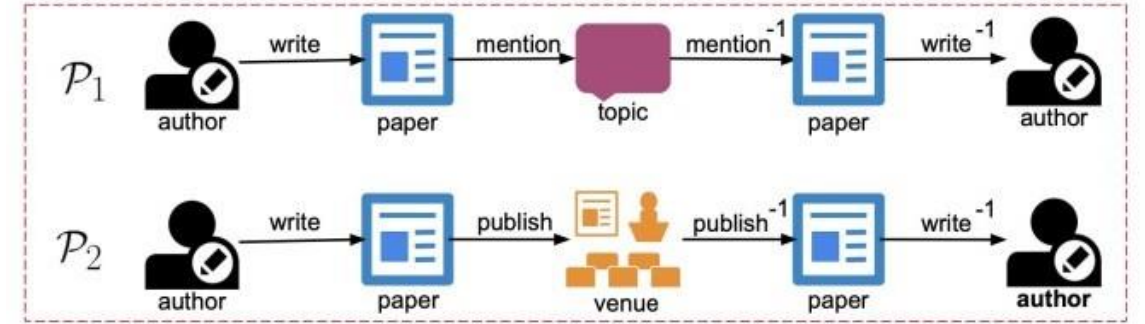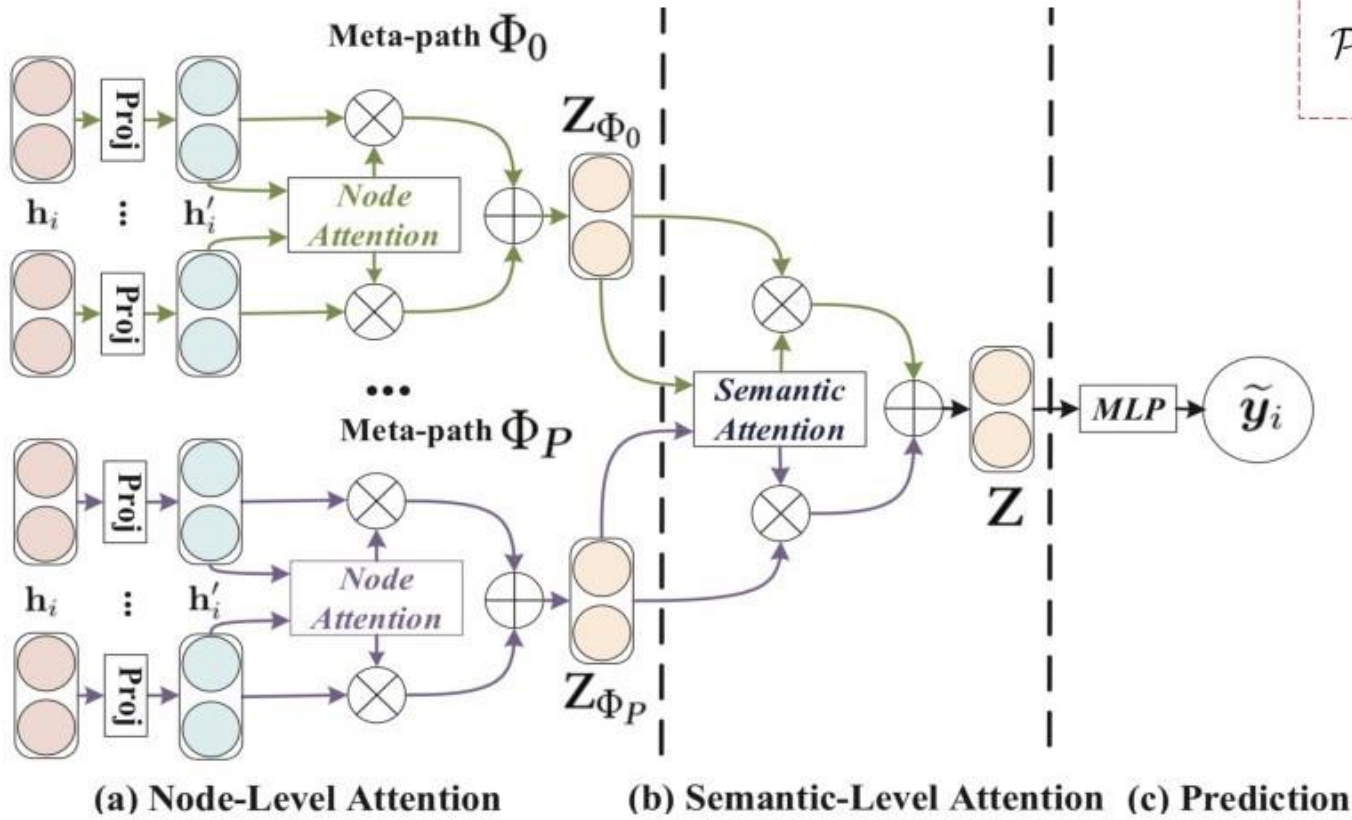$$W_r^{(l)} = \sum_{b=1}^{B} a_{rb}^{(l)} V_b^{(l)}$$

  - Block decomposition

$$W_r^{(l)} = \bigoplus_{b=1}^{B} Q_{br}^{(l)}.$$
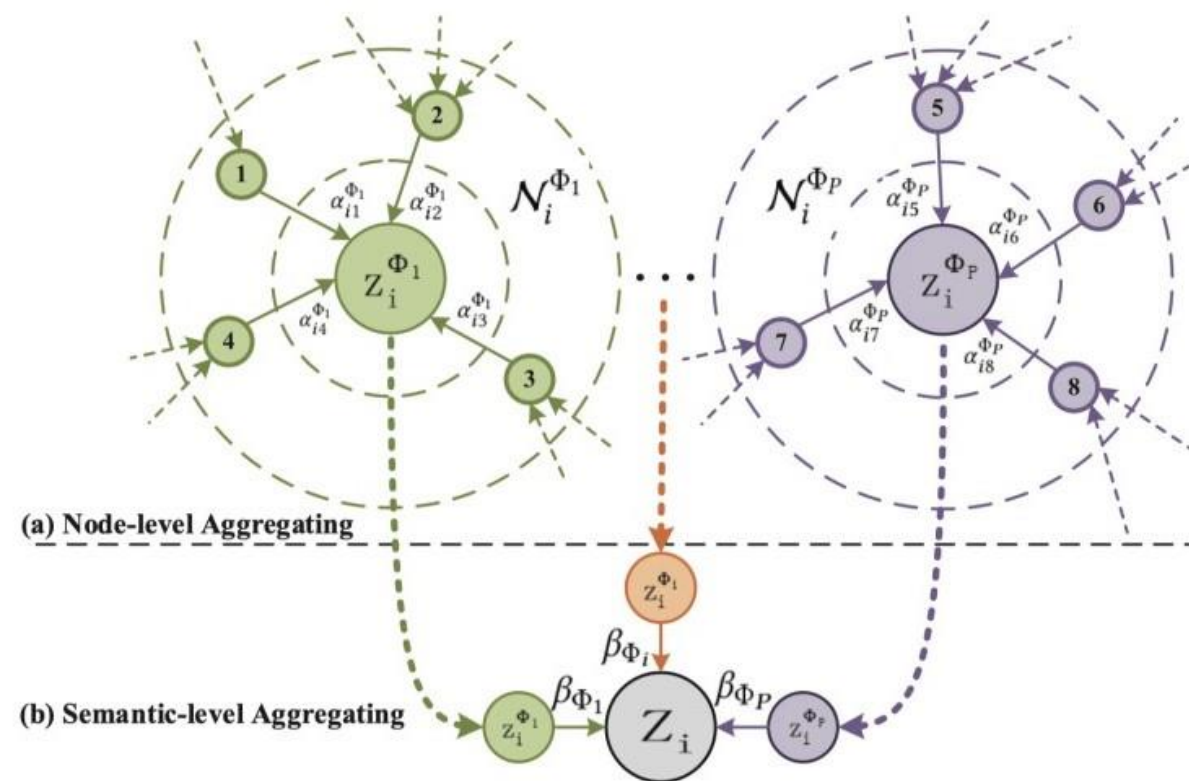
# HAN

- Meta path



Meta paths among author nodes



(a) Node-Level Attention  (b) Semantic-Level Attention  (c) Prediction

- Node Level Attention

$$\alpha_{ij}^{\Phi} = softmax_j(e_{ij}^{\Phi}) = \frac{\exp(\sigma(\mathbf{a}_{\Phi}^{\mathrm{T}} \cdot [\mathbf{h}_i' \| \mathbf{h}_j']))}{\sum_{k \in \mathcal{N}_i^{\Phi}} \exp(\sigma(\mathbf{a}_{\Phi}^{\mathrm{T}} \cdot [\mathbf{h}_i' \| \mathbf{h}_k']))}$$

$$\mathbf{z}_i^{\Phi} = \overset{K}{\underset{k=1}{\|}} \sigma\left( \sum_{j \in \mathcal{N}_i^{\Phi}} \alpha_{ij}^{\Phi} \cdot \mathbf{h}_j' \right).$$
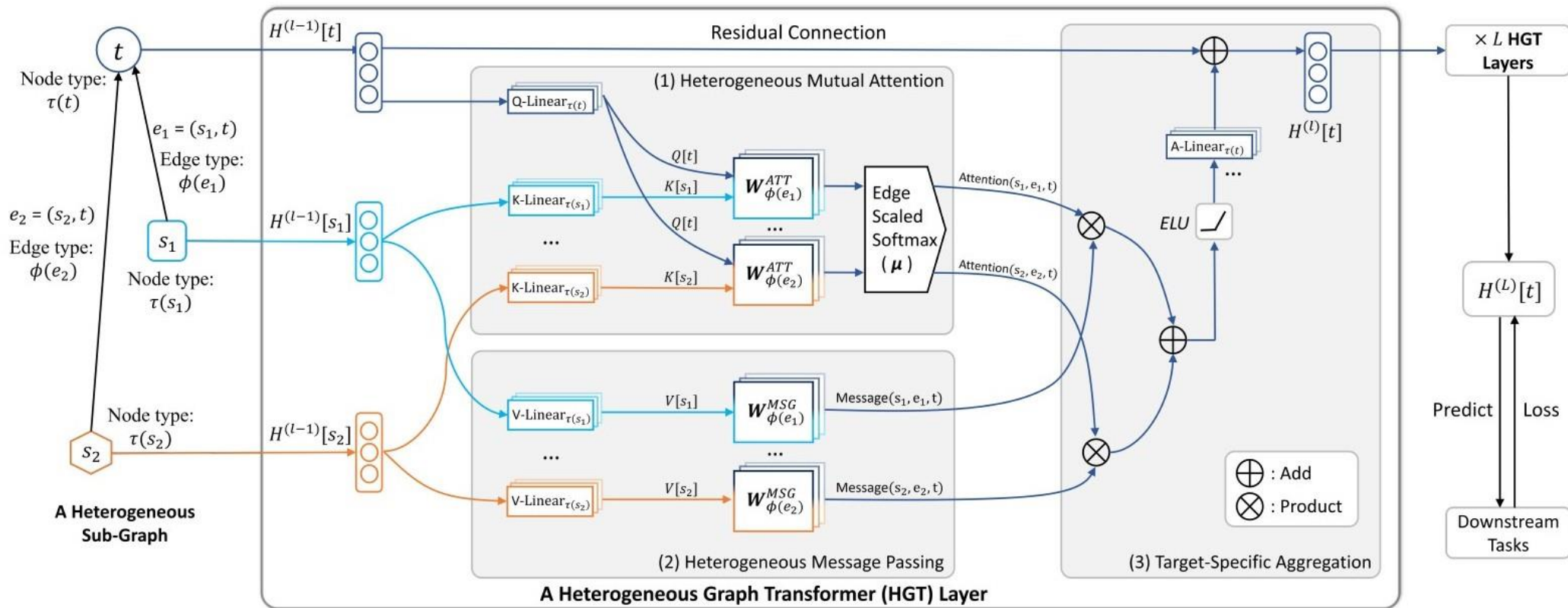
# HAN

- Semantic Level Attention



$$w_{\Phi_p} = \frac{1}{|\mathcal{V}|} \sum_{i \in \mathcal{V}} \mathbf{q}^{\mathrm{T}} \cdot \tanh(\mathbf{W} \cdot \mathbf{z}_i^{\Phi_p} + \mathbf{b})$$

$$\beta_{\Phi_p} = \frac{\exp(w_{\Phi_p})}{\sum_{p=1}^{P} \exp(w_{\Phi_p})}$$

$$\mathbf{Z} = \sum_{p=1}^{P} \beta_{\Phi_p} \cdot \mathbf{Z}_{\Phi_p}$$

(a) Node-level Aggregating

(b) Semantic-level Aggregating

# HGT



$$H^l[t] \leftarrow \underset{\forall s \in N(t), \forall e \in E(s,t)}{\textbf{Aggregate}} \left( \textbf{Attention}(s,t) \cdot \textbf{Message}(s) \right)$$

$$\textbf{Attention}_{HGT}(s,e,t) = \underset{\forall s \in N(t)}{\text{Softmax}} \left( \underset{i \in [1,h]}{\|} ATT\text{-}head^i(s,e,t) \right) \quad (3)$$

$$ATT\text{-}head^i(s,e,t) = \left( K^i(s) \, W^{ATT}_{\phi(e)} \, Q^i(t)^T \right) \cdot \frac{\mu_{\langle \tau(s),\phi(e),\tau(t) \rangle}}{\sqrt{d}}$$

$$\textbf{Message}_{HGT}(s,e,t) = \underset{i \in [1,h]}{\|} MSG\text{-}head^i(s,e,t)$$

$$K^i(s) = \text{K-Linear}^i_{\tau(s)} \left( H^{(l-1)}[s] \right)$$

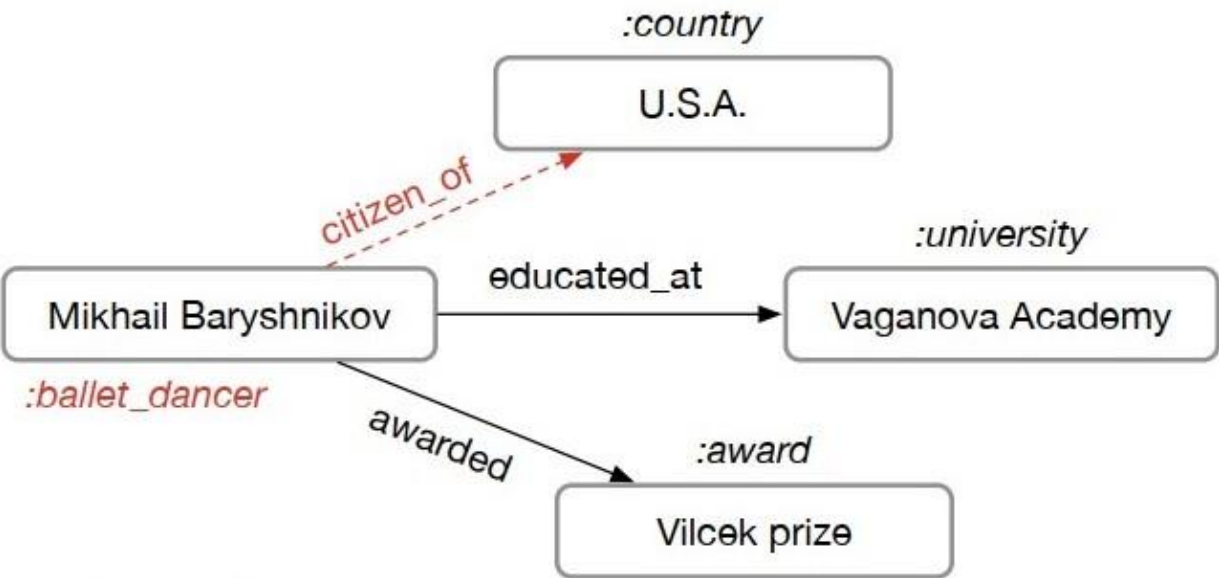$$MSG\text{-}head^i(s,e,t) = \text{M-Linear}^i_{\tau(s)} \left( H^{(l-1)}[s] \right) W^{MSG}_{\phi(e)}$$

$$Q^i(t) = \text{Q-Linear}^i_{\tau(t)} \left( H^{(l-1)}[t] \right)$$

# GNN for KG Embedding

$$h_s = \sigma\left(\sum_{r\in\mathcal{R}}\sum_{o\in\mathcal{N}_r(s)} f\left(s,r,o\right)\right)$$



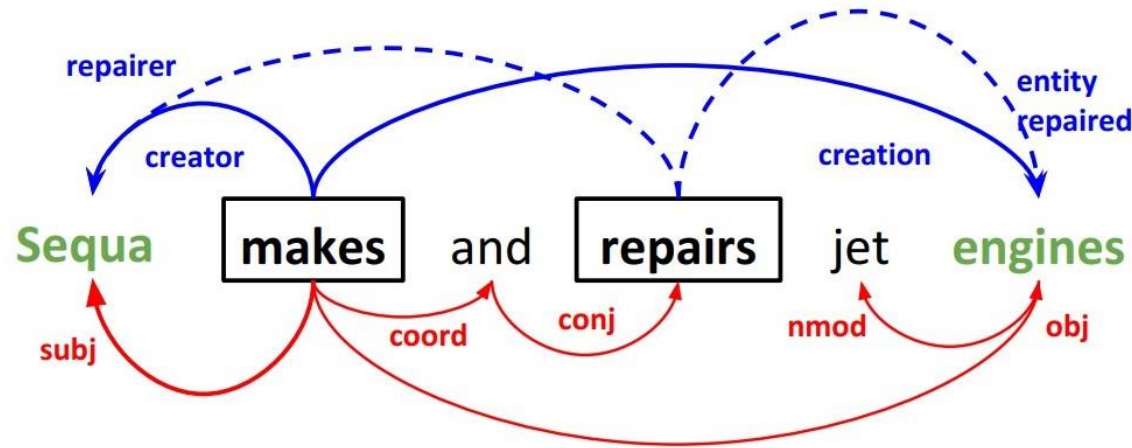| Method | $f\left(s,r,o\right)$ |
|---|---|
| R-GCN | $W_r h_o$ |
| SACN | $\alpha_r W h_o$ |
| KBGAT | $\alpha_{sro} W h_{sro}$ |
| CompGCN | $W_{\lambda(r)}\phi(s,r)$ |

# GNN for KG Embedding: Encoder-Decoder for LP

# Application: Syntactic GCN for SRL

- Task: discover who did what to whom.
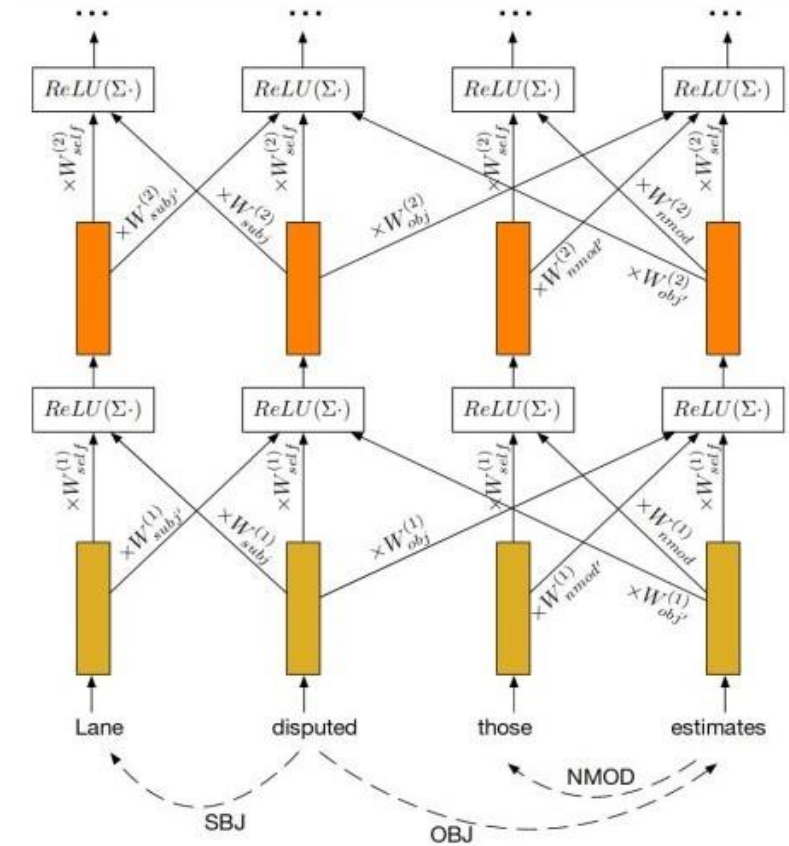- Syntax mirrors semantics
- Exploit syntax using convolution



$$h_v = ReLU\left(\sum_{u \in \mathcal{N}(v)} g_{u,v}\left(W_{d(u,v)} h_u + b_{l(u,v)}\right)\right)$$

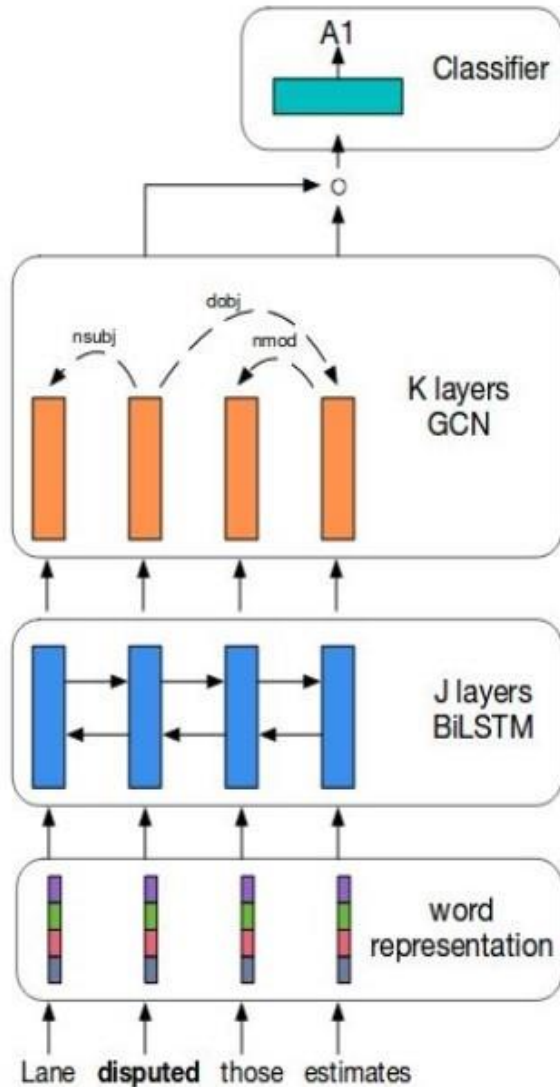word emb of $v$      weight of direction      bias of label + direction

**edge-wise gating**

$$g_{u,v} = \sigma\left(\hat{w}_{d(u,v)} h_u + \hat{b}_{l(u,v)}\right)$$

# Application: Syntactic GCN for SRL



- Trained with cross-entropy loss

**F1 on CoNLL-2009**

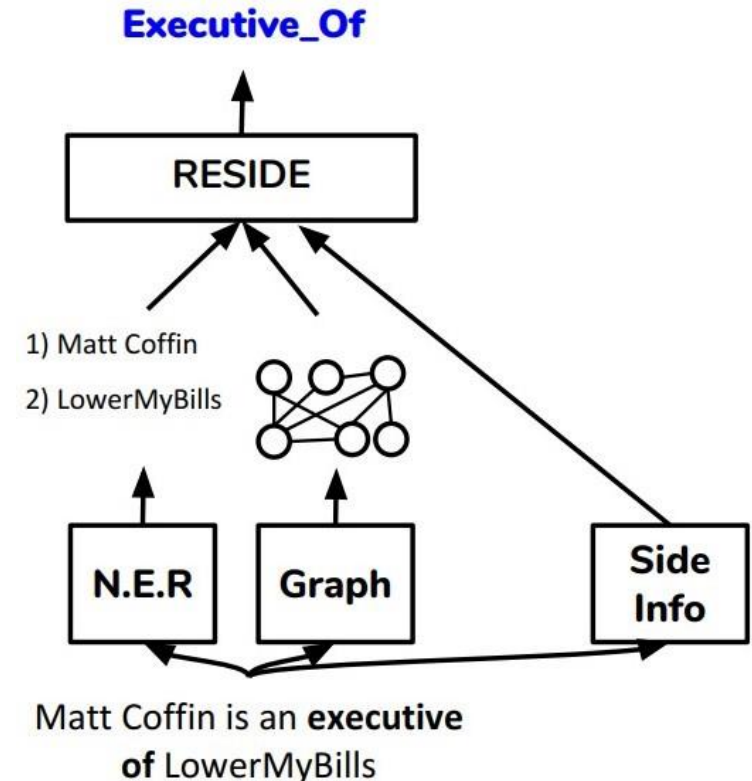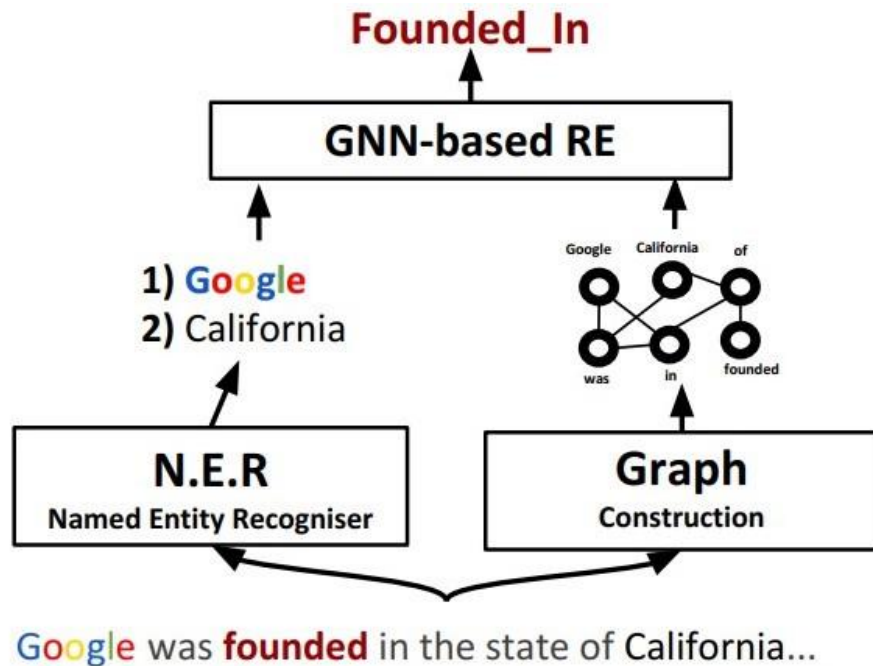| BiLSTM | 82.7 |
|---|---|
| BiLSTM + GCN | **83.3** |

- GCN integrates syntax, context
- GCN, LSTM complement each other

# Application: GNNs for Relation Extraction

- Identify relation between entities

# GNNs for Relation Extraction: RE-SIDE



- Side information Acquisition
  - Relation Alias side information
    Open IE, PPDB -> GloVe Embeddings ->  find the
    closest relation for each phase (h^{rel})

  - Entity Type side information
    E.g. Paris: government, location -> average type
    embeddings (h^{type})

- Side info improves performance!

# Application: QA-GNN

If it is not used for **hair**, a **round brush** is an example of what?
A. hair brush  B. bathroom  C. **art supplies***
D. shower  E. hair salon

**QA context**

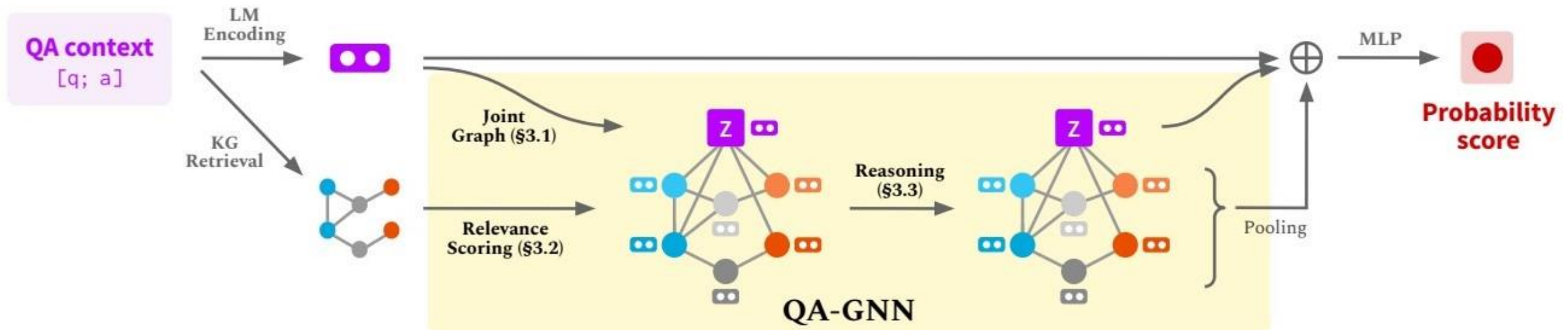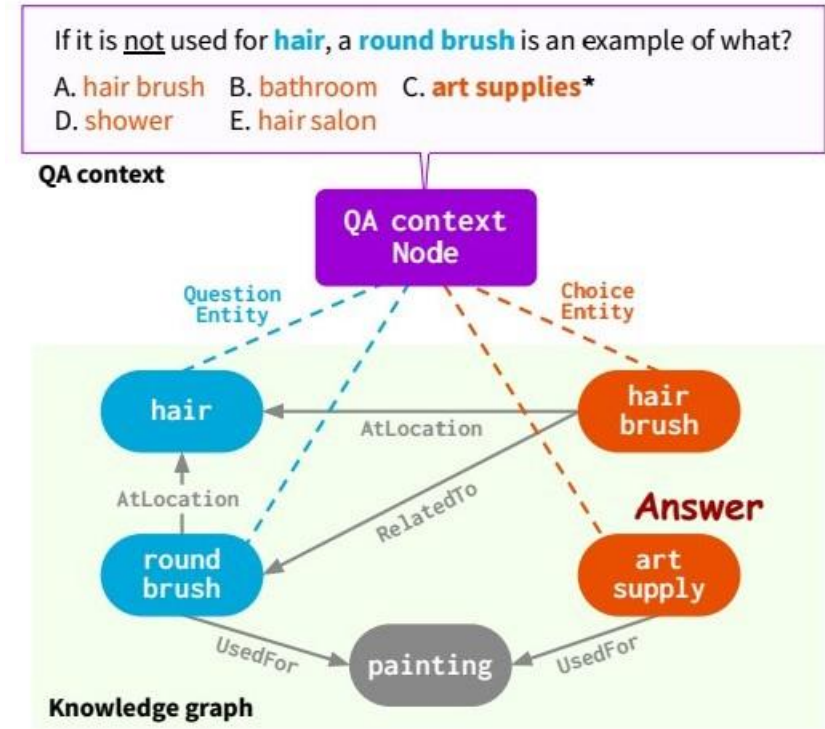- Identify relevant knowledge from KG -> **relevance scoring**
- Performa **joint reasoning** over the QA context and KG



QA context Node

Question Entity          Choice Entity

hair          AtLocation          hair brush

AtLocation          RelatedTo          Answer

round brush          art supply

UsedFor          painting          UsedFor

**Knowledge graph**



**QA context** [q; a]

LM Encoding

KG Retrieval

Joint Graph (§3.1)

Relevance Scoring (§3.2)

Z

Reasoning (§3.3)

Z

Pooling

MLP

**Probability score**

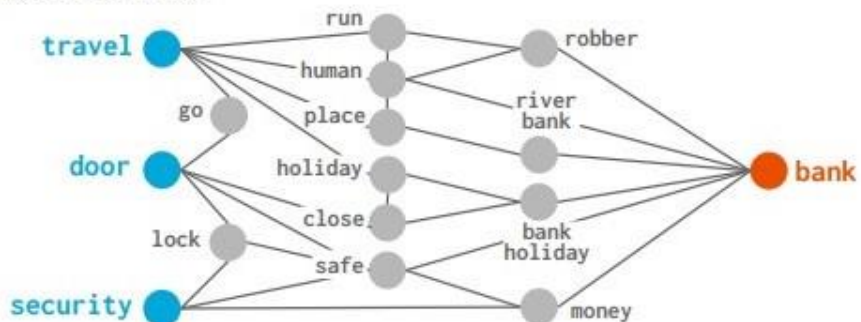**QA-GNN**

# QA-GNN: relevance scoring



**QA Context**

A **revolving door** is convenient for **two direction travel**, but also serves as a **security measure** at what?

A. **bank***    B. library    C. department store
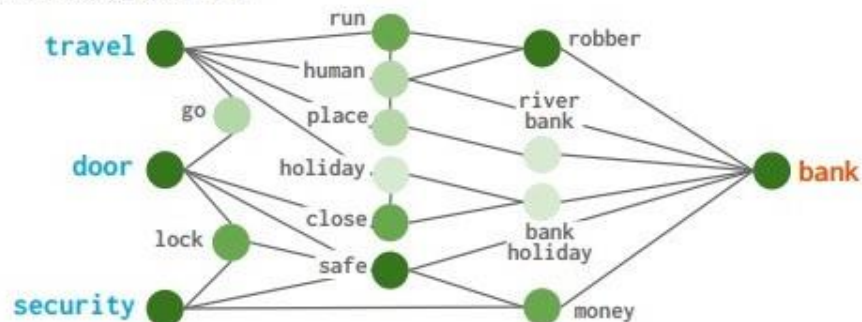D. mall    E. new york

**Language Model**

Relevance ( entity | QA context )

entity

**Retrieved KG**

travel, run, robber, human, go, place, river bank, door, holiday, close, bank, lock, bank holiday, safe, security, money, bank

Some entities are more relevant than others given the context.

**KG node scored**

travel, run, robber, human, go, place, river bank, door, holiday, close, bank, lock, bank holiday, safe, security, money, bank

Entity relevance estimated. **Darker** color indicates higher score.

$$\rho_v = f_{\text{head}}(f_{\text{enc}}([\text{text}(z); \text{text}(v)])),$$
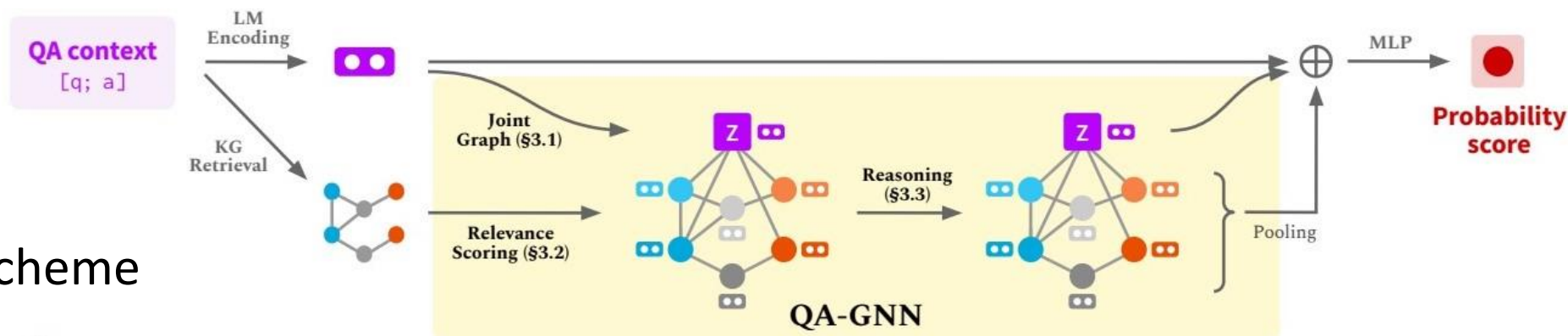
# QA-GNN: GNN architecture



- Message passing scheme

$$h_t^{(\ell+1)} = f_n \left( \sum_{s \in \mathcal{N}_t \cup \{t\}} \alpha_{st} m_{st} \right) + h_t^{(\ell)}$$

- Message

$$m_{st} = f_m(h_s^{(\ell)}, u_s, r_{st})$$

$$u_t = f_u(u_t), \quad r_{st} = f_r(e_{st}, u_s, u_t)$$

- Attention weight

$$\rho_t = f_\rho(\rho_t)$$

$$q_s = f_q(h_s^{(\ell)}, u_s, \rho_s),$$

$$k_t = f_k(h_t^{(\ell)}, u_t, \rho_t, r_{st})$$

$$\alpha_{st} = \frac{\exp(\gamma_{st})}{\sum_{t' \in \mathcal{N}_s \cup \{s\}} \exp(\gamma_{st'})}, \quad \gamma_{st} = \frac{q_s^\top k_t}{\sqrt{D}}$$

# Thx for Attention

Zhaoxuan Tan

Xi'an Jiaotong University

March 5, 2022