

# DSC250: Advanced Data Mining

## Knowledge Graphs (KGs)


**Zhiting Hu**

Lecture 14, Feb 20, 2025

**UC San Diego**

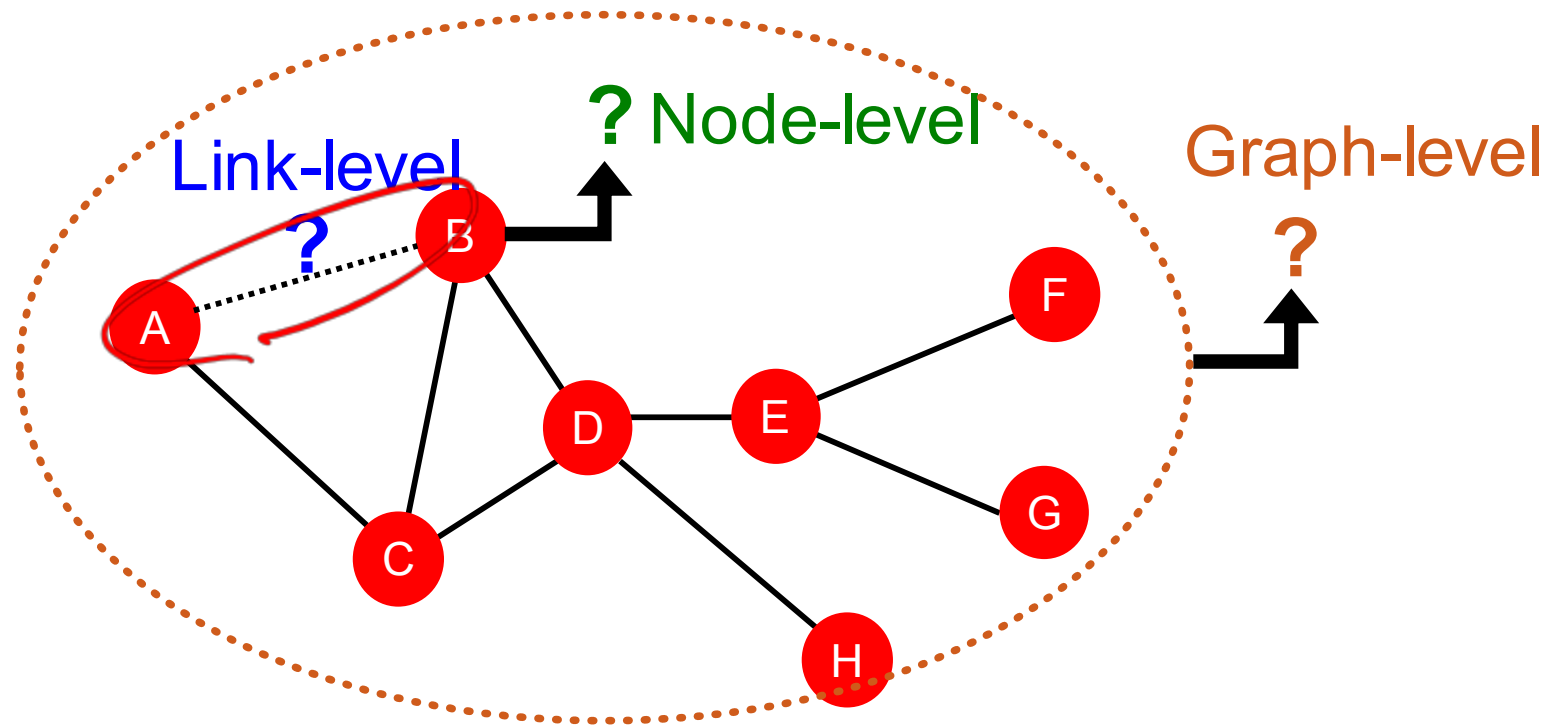
**HALICIOĞLU DATA SCIENCE INSTITUTE**

# Outline

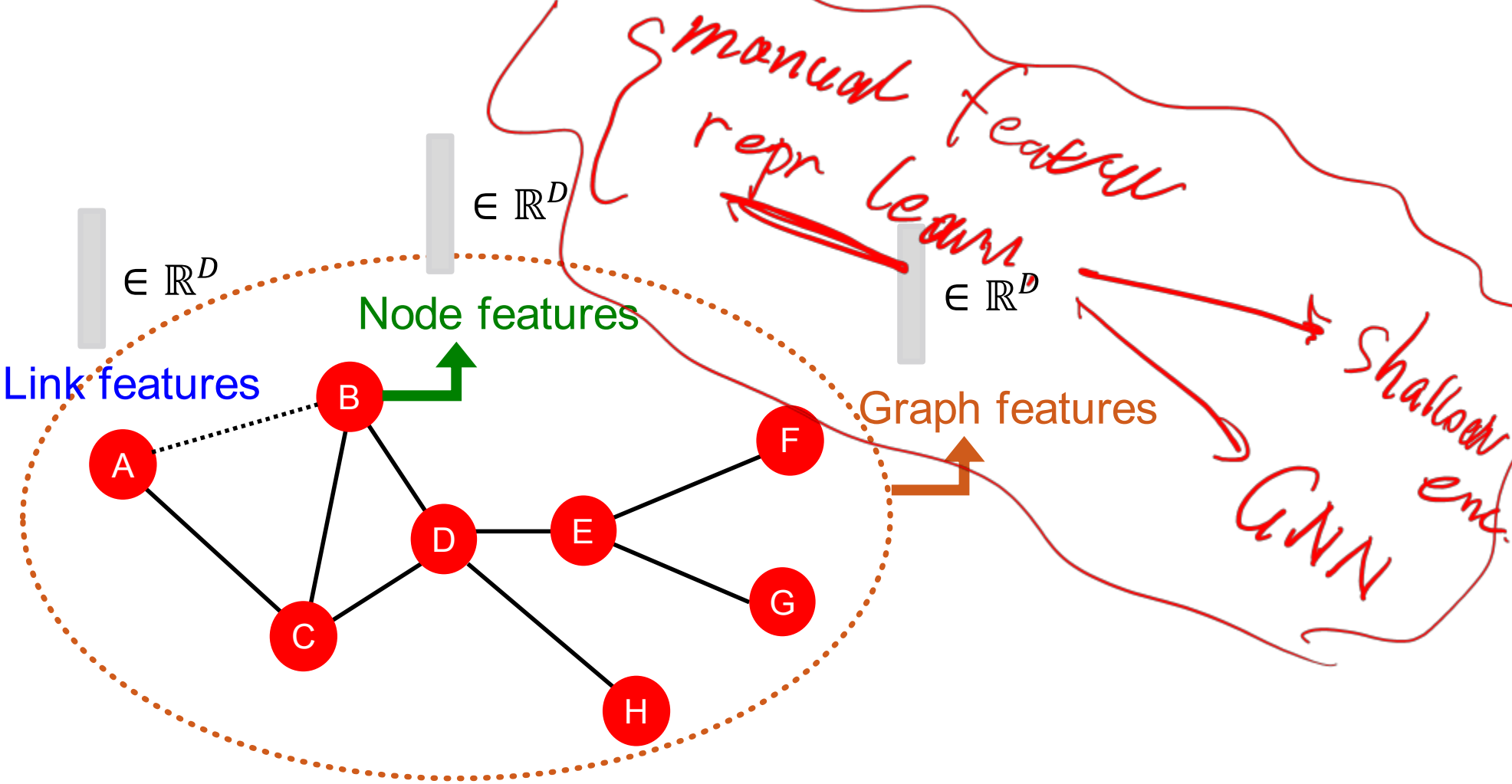
- Knowledge graphs
  - Presentation
    - Lingyi Yang, Jiayue Yuan: "Generative Agents: Interactive Simulacra of Human Behavior"
    - Yuyuan Wu, Zijie Feng: "Data-efficient Fine-tuning for LLM-based Recommendation"
    - Ivy Nguyen, Joshua Chuang: "Prediction of COVID-19 cases by multifactor driven long short-term memory (LSTM) model"
    - Shuyu Wang, Caroline Zhang: "Yo'LLaVA: Your Personalized Language and Vision Assistant"
- 

# Recap: Tasks on Graph

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



# Recap: Getting Features for Nodes/Links/Graphs

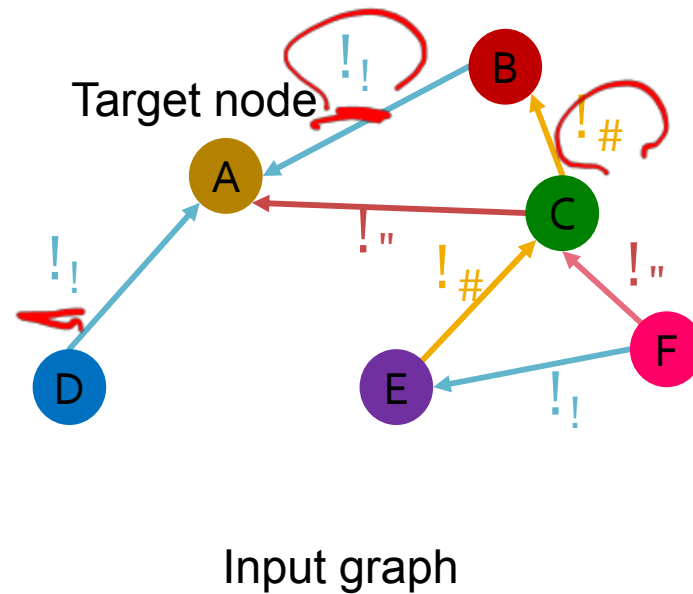


# Outline

- Overview
- Knowledge Graph Completion (Link Prediction)
- Reasoning on Knowledge Graphs

# Heterogeneous Graphs

- ! **Heterogeneous graphs:** a graph with **multiple relation types**



# Heterogeneous Graphs

8 possible relation types!

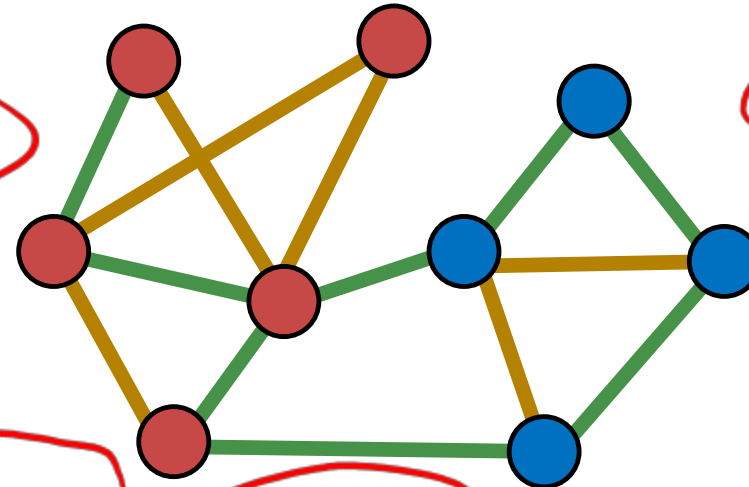


(Paper, Cite, Paper)

(Paper, Like, Paper)

(Paper, Cite, Author)

(Paper, Like, Author)



(Author, Cite, Author)

(Author, Like, Author)

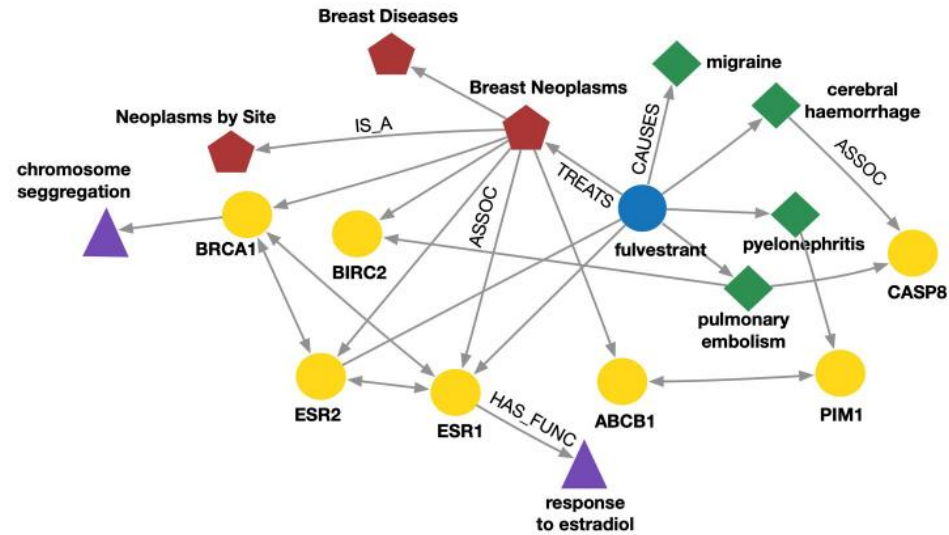
(Author, Cite, Paper)

(Author, Like, Paper)

**Relation types:** (node\_start, edge, node\_end)

- ! We use **relation type to describe an edge** (as opposed to edge type)
- ! Relation type better captures the interaction between nodes and edges

# Heterogeneous Graphs



*prevent*  
*disease*  
*drug*  
*treat*

## Biomedical Knowledge Graphs

**Example node:** Migraine

**Example relation:** (fulvestrant, Treats, Breast Neoplasms)

**Example node type:** Protein

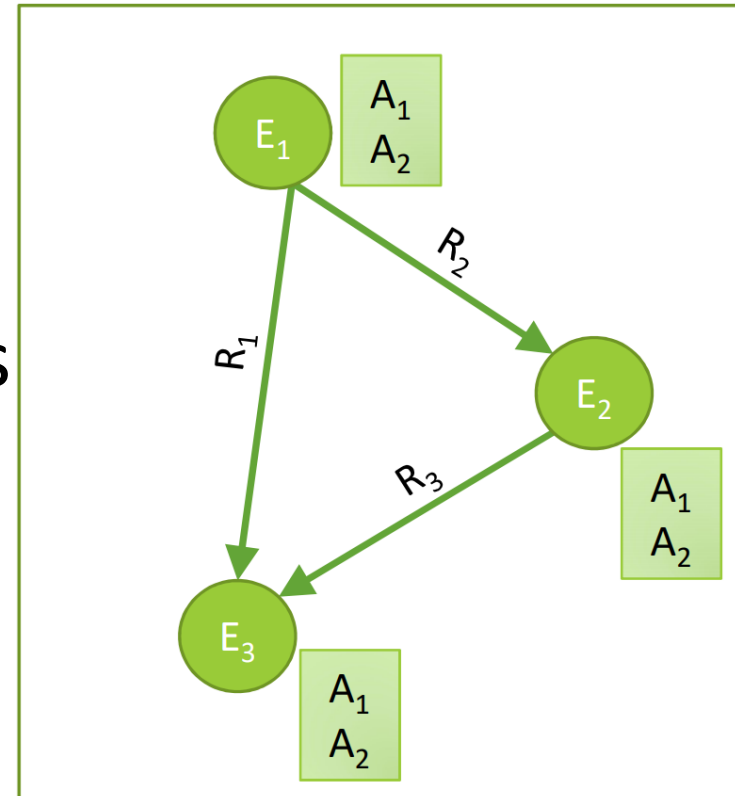
**Example edge type:** Causes



# Knowledge Graph

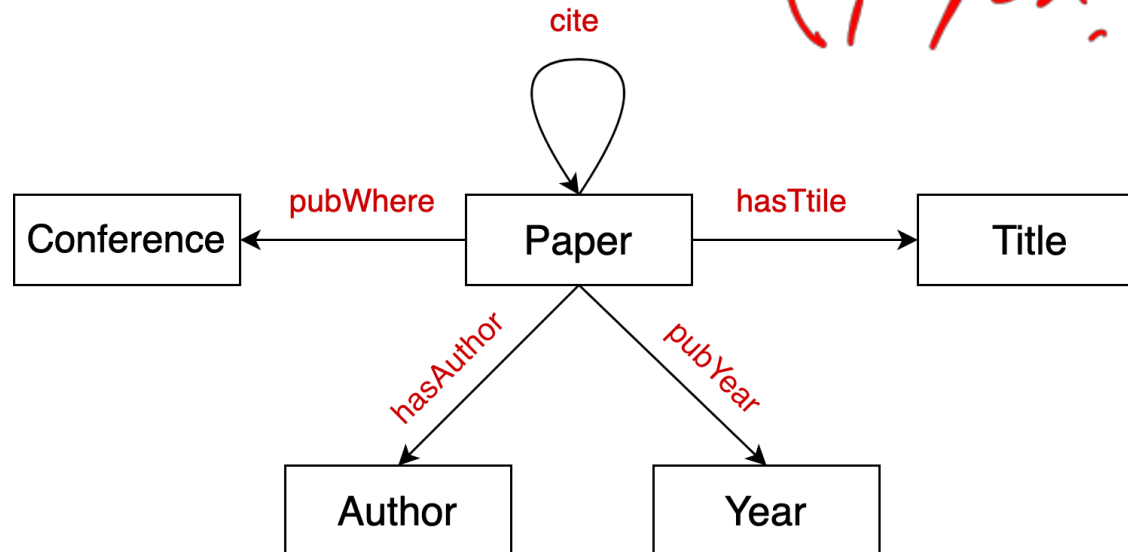
## Knowledge in graph form:

- § Capture entities, types, and relationships
- ! Nodes are **entities**
- ! Nodes are labeled with their **types**
- ! Edges between two nodes capture **relationships** between entities
- ! **KG is an example of a heterogeneous graph**



# Example: Bibliographic Networks

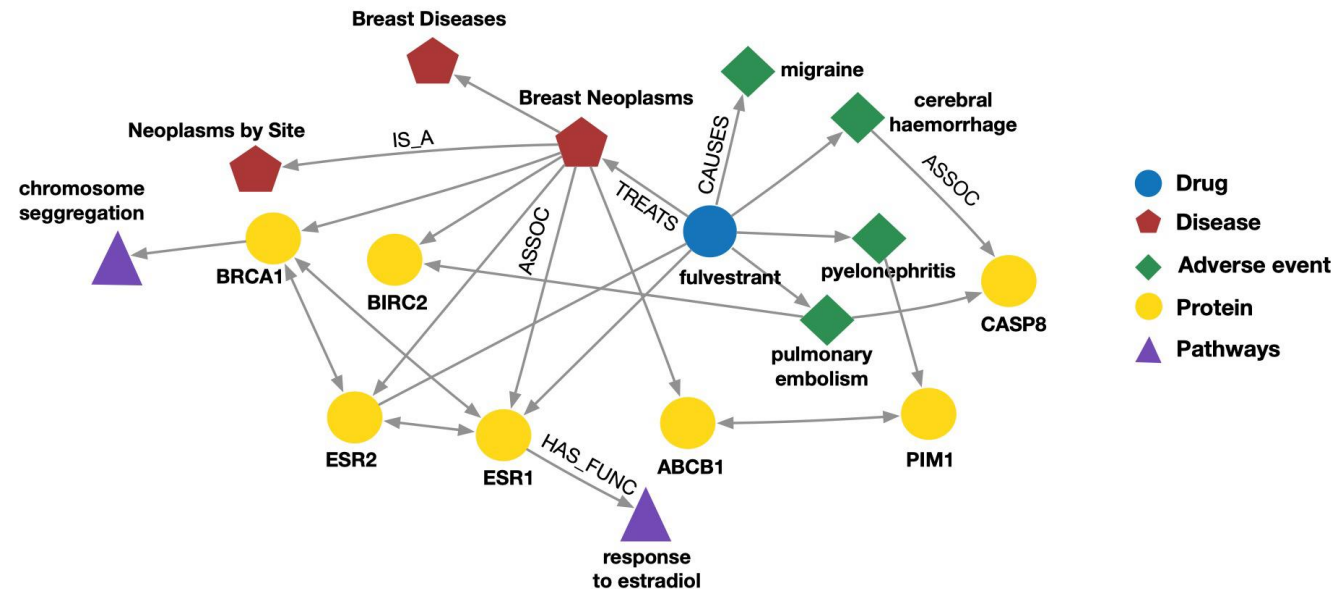
- Node types: paper, title, author, conference, year
- Relation types: pubWhere, pubYear, hasTitle, hasAuthor, cite



*Handwritten notes:*  
(paper, where, conf)  
node\_start: Paper  
node\_end: title

# Example: Bio Knowledge Graphs

- Node types: drug, disease, adverse event, protein, pathways
- Relation types: has\_func, causes, assoc, treats, is\_a



# KGs in Practice

## Examples of knowledge graphs

- ! Google Knowledge Graph
- ! Amazon Product Graph
- ! Facebook Graph API
- ! IBM Watson
- ! Microsoft Satori
- ! Project Hanover/Literome
- ! LinkedIn Knowledge Graph
- ! Yandex Object Answer

# Applications of KGs

## i Serving information:

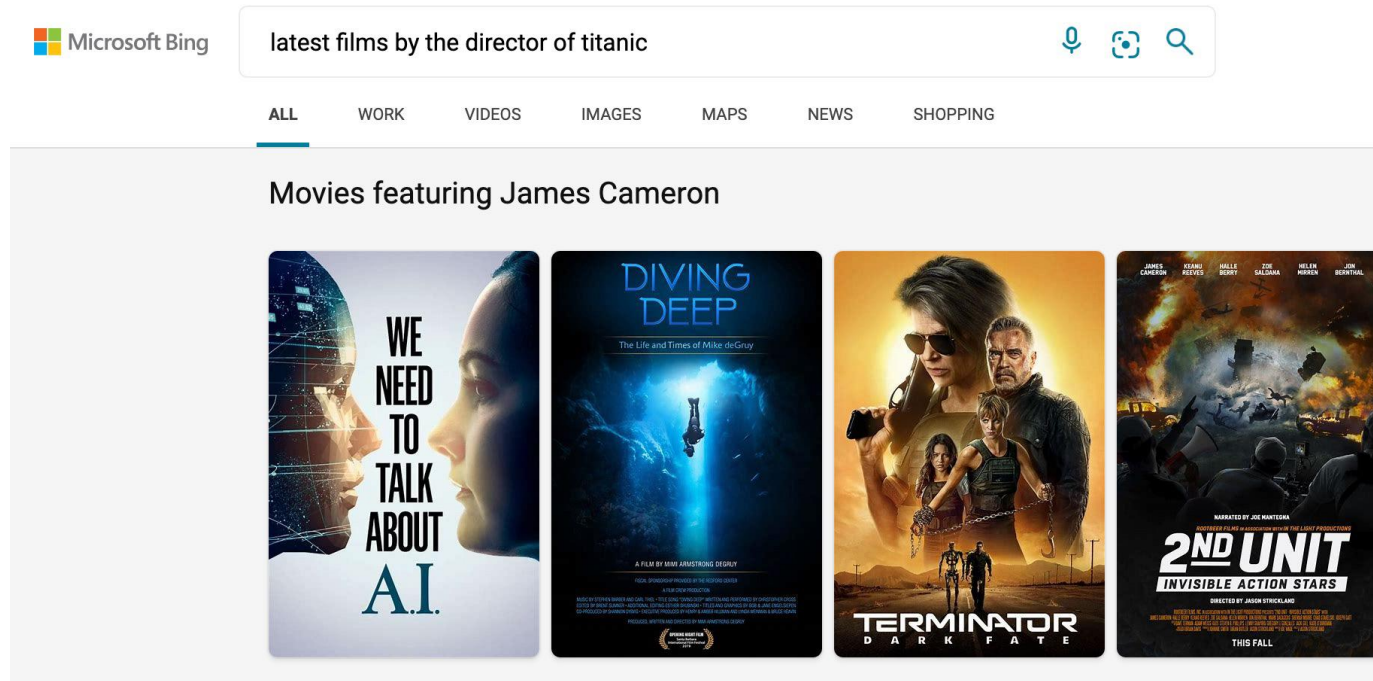
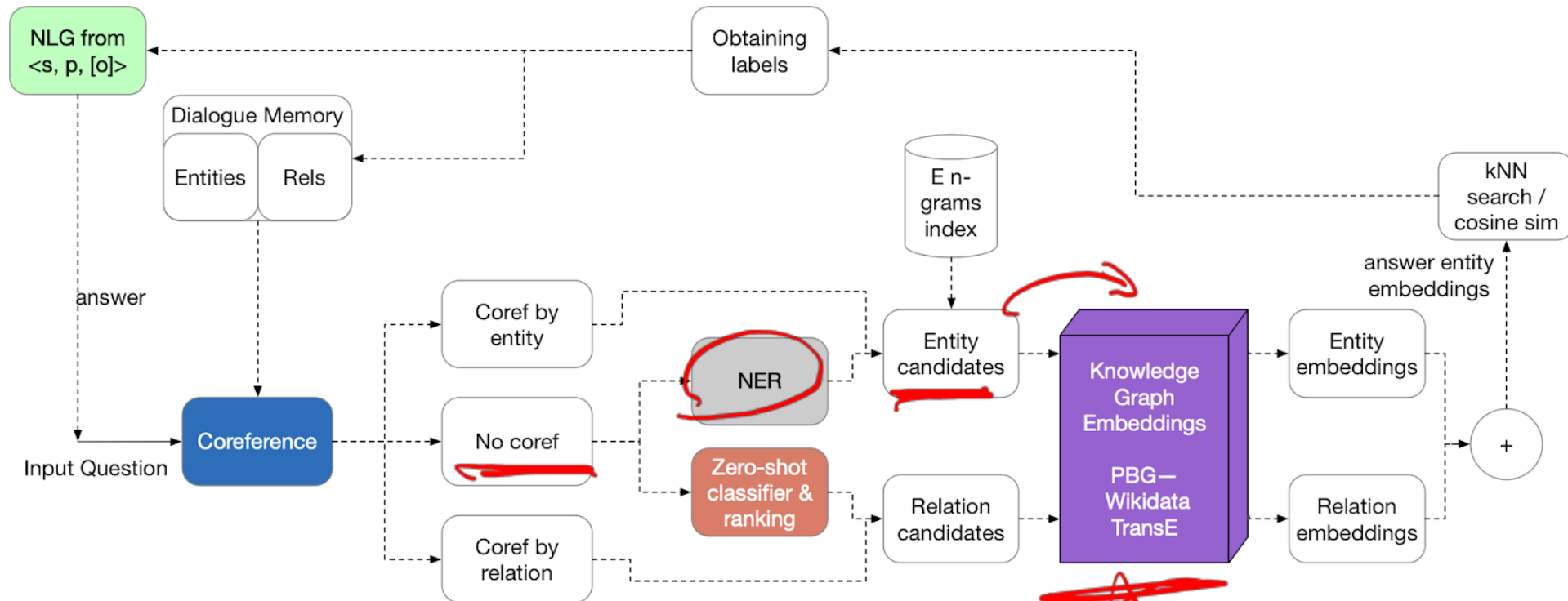


Image credit: Bing

# Applications of KGs

*KG* *RAG* *LLM*

## Question answering and conversation agents



*Core NLP*

Image credit: [Medium](#)

# KG Datasets

ConceptNet

## Publicly available KGs:

§ FreeBase, Wikidata, Dbpedia, YAGO, NELL, etc.

## Common characteristics:

§ **Massive**: Millions of nodes and edges

§ **Incomplete**: Many true edges are missing

Bill Gates funded MS  
Never-ending Language Learning



# KG Datasets

- Publicly available KGs:

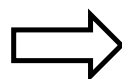
- § FreeBase, Wikidata, Dbpedia, YAGO, NELL, etc.

- Common characteristics:

- § **Massive**: Millions of nodes and edges

- § **Incomplete**: Many true edges are missing

Given a massive KG,  
enumerating all the  
possible facts is  
intractable!



Can we predict plausible  
BUT missing links?





# Example: Freebase



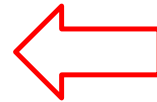
## i Freebase

§ ~80 million entities

§ ~38K relation types

§ ~3 billion facts/triples

93.8% of persons from Freebase  
have no place of birth and 78.5%  
have no nationality!



## i Datasets: FB15k/FB15k-237

§ A **complete** subset of Freebase, used by researchers to learn KG models

Dataset	Entities	Relations	Total Edges
FB15k	14,951	1,345	592,213
FB15k-237	14,505	237	310,079

[1] Paulheim, Heiko. "Knowledge graph refinement: A survey of approaches and evaluation methods." *Semantic web* 8.3 (2017): 489-508.

[2] Min, Bonan, et al. "Distant supervision for relation extraction with an incomplete knowledge base." *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2013.

# Outline

- Overview
- **Knowledge Graph Completion (Link Prediction)**
- Reasoning on Knowledge Graphs



# KG Representation

- Edges in KG are represented as **triples**  $(h, r, t)$   
§ **head** ( $h$ ) has **relation** ( $r$ ) with **tail** ( $t$ )

# KG Representation

- i Edges in KG are represented as **triples**  $(h, r, t)$ 
  - § head ( $h$ ) has relation ( $r$ ) with tail ( $t$ )
- i **Key Idea:**
  - § Model entities and relations in embedding space  $\mathbb{R}^d$ 
    - § Associate entities and relations with **shallow embeddings**
    - § **Note we do not learn a GNN here!**

# KG Representation

i Edges in KG are represented as **triples**  $(h, \$, \%)$

§ head ( $h$ ) has relation ( $\$$ ) with tail ( $\%$ )

i **Key Idea:**

§ Model entities and relations in embedding space  $\mathbb{R}^d$

§ Associate entities and relations with **shallow embeddings**

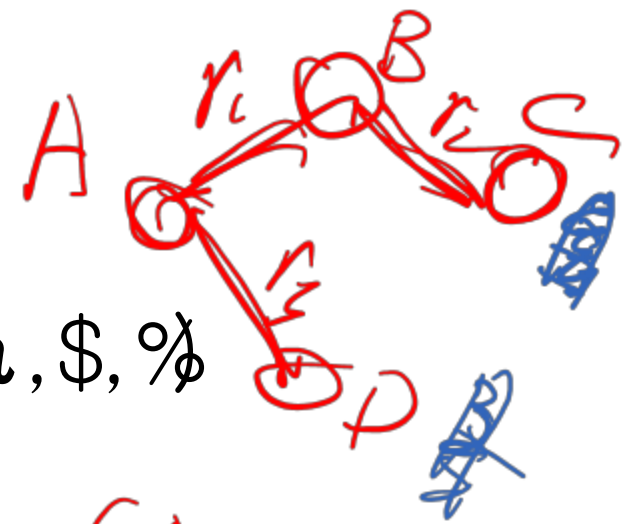
§ **Note we do not learn a GNN here!**

§ Given a triple  $(h, \$, \%)$ , the goal is that the **embedding of  $(h, \$)$  should be close to the embedding of  $\%$**

§ How to embed  $(h, \$)$ ?

§ How to define score  $s(h, \%)$ ?

§ Score  $s$  is high if  $(h, \$, \%)$  exists, else  $s$  is low



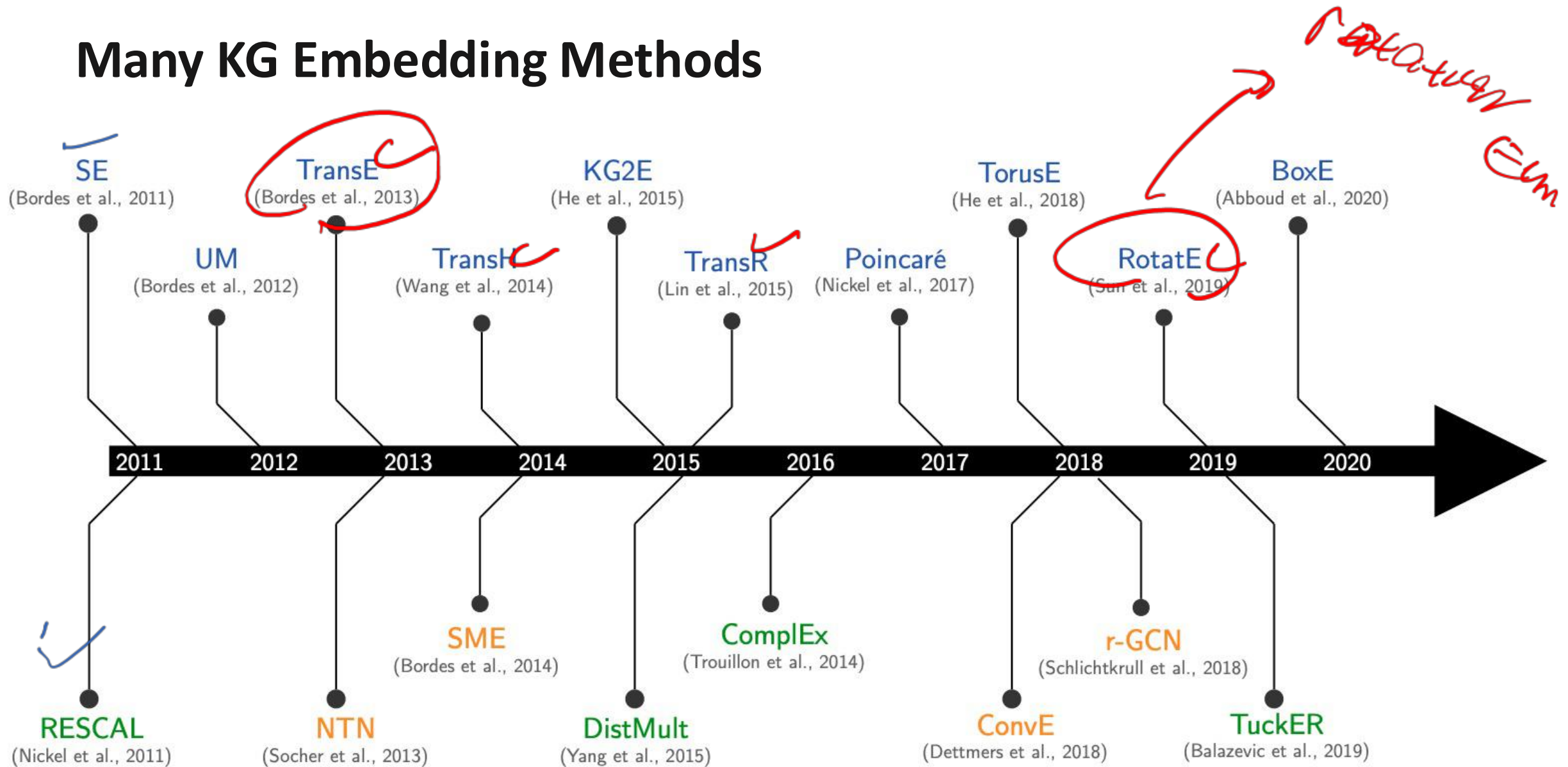
$(A, r_1, B)$

$(B, r_2, C)$

$(A, r_3, D)$

$(A, B, ?)$

# Many KG Embedding Methods



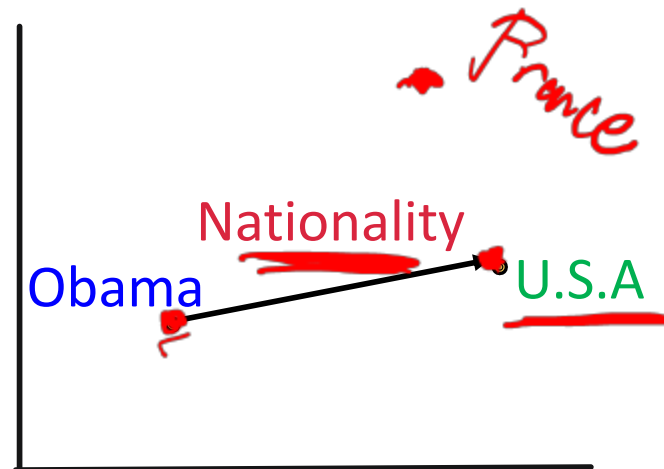
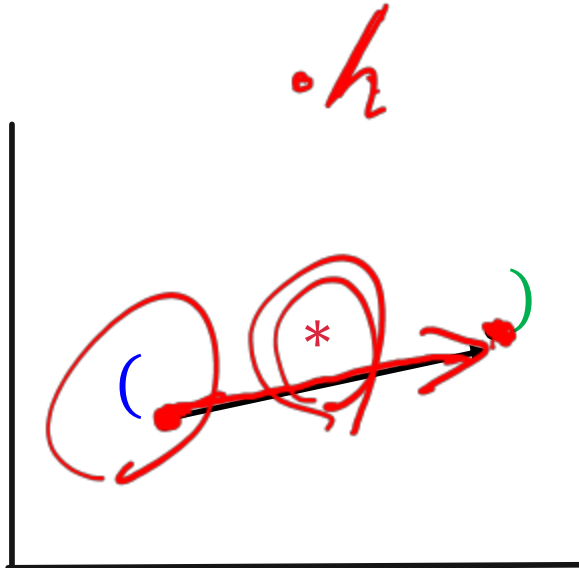
# TransE for KG Completion

## i Intuition: Translation

For a triple  $(h, r, t)$ , let  $\mathbf{h}, \mathbf{r}, \mathbf{t} \in \mathbb{R}^d$  be embedding vectors.

embedding vectors  
will appear in  
boldface

i **TransE:  $\mathbf{h} + \mathbf{r} \approx \mathbf{t}$**  if the given link exists else  $\mathbf{h} + \mathbf{r} \neq \mathbf{t}$





# TransE for KG Completion

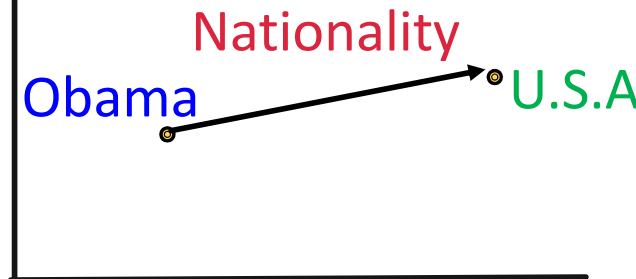
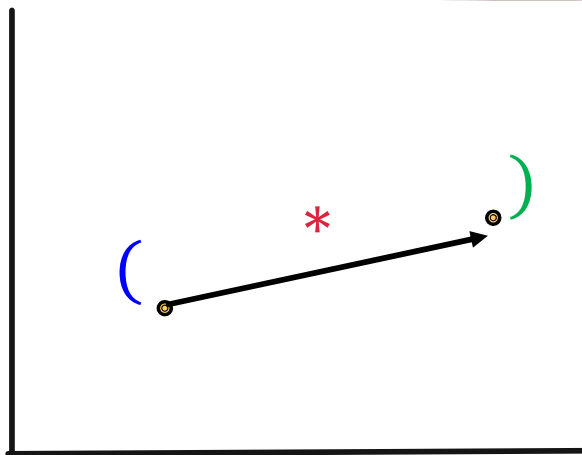
## i Intuition: Translation

For a triple  $(h, r, t)$ , let  $\mathbf{h}, \mathbf{r}, \mathbf{t} \in \mathbb{R}^d$  be embedding vectors.

embedding vectors  
will appear in  
boldface

i **TransE**:  $\mathbf{h} + \mathbf{r} \approx \mathbf{t}$  if the given link exists else  $\mathbf{h} + \mathbf{r} \neq \mathbf{t}$

Entity scoring function:  $f_r(h, t) = -\|\mathbf{h} + \mathbf{r} - \mathbf{t}\|$



$$f_r(h, t_1)$$

$$f_r(h, t_2)$$

# Connectivity Patterns in KG

- i Relations in a heterogeneous KG have different properties:

§ Example:

§ Symmetry: If the edge  $(h, \text{"Roommate"}, \phi)$  exists in KG, then the edge  $(\phi, \text{"Roommate"}, h)$  should also exist.

§ Inverse relation: If the edge  $(h, \text{"Advisor"}, \phi)$  exists in KG, then the edge  $(\phi, \text{"Advisee"}, h)$  should also exist.

## Connectivity Patterns in KG

- i **Relations in a heterogeneous KG have different properties:**

§ Example:

§ **Symmetry:** If the edge  $(h, \text{"Roommate"}, \phi)$  exists in KG, then the edge  $(\phi, \text{"Roommate"}, h)$  should also exist.

§ **Inverse relation:** If the edge  $(h, \text{"Advisor"}, \phi)$  exists in KG, then the edge  $(\phi, \text{"Advisee"}, h)$  should also exist.

- i **Can we categorize these relation patterns?**
- i **Are KG embedding methods (e.g., TransE) expressive enough to model these patterns?**

# Four Relationship Patterns

## i **Symmetric (Antisymmetric) Relations:**

$$r(h, t) \Rightarrow r(t, h) \quad (\underline{r(h, t)} \Rightarrow \underline{\neg r(t, h)}) \quad \forall h, t$$

### § **Example:**

§ Symmetric: Family, Roommate

§ Antisymmetric: Hypernym (a word with a broader meaning: poodle vs. dog)

## i **Inverse Relations:**

$$r_l(h, t) \Rightarrow r_j(t, h)$$

§ **Example** : (Advisor, Advisee)

## i **Composition (Transitive) Relations:**

$$\underline{r_j(x, y)} \wedge \underline{r_l(y, z)} \Rightarrow \underline{r_*(x, z)} \quad \forall x, y, z$$

§ **Example:** My mother's husband is my father.

## i **1-to-N relations:**

$$\underline{r(h, t_j), r(h, t_l), \dots, r(h, t_+)} \text{ are all True.}$$

§ **Example:** # is "StudentsOf"

# Antisymmetric Relations in TransE

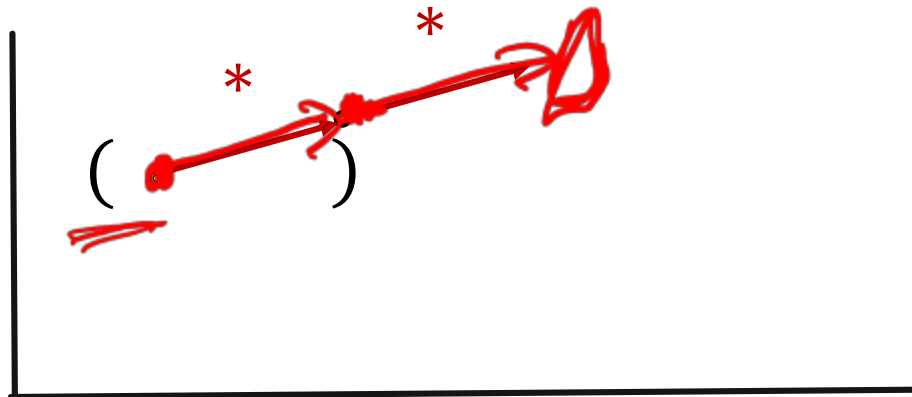
## i Antisymmetric Relations:

$$\$(h, r) \Rightarrow \neg \$(r, h) \square \forall h, r$$

§ **Example:** Hypernym (a word with a broader meaning: poodle vs. dog)

## i **TransE** can model antisymmetric relations **ü**

§ ( + \* = ), but ) + \* ≠ (



# Inverse Relations in TransE

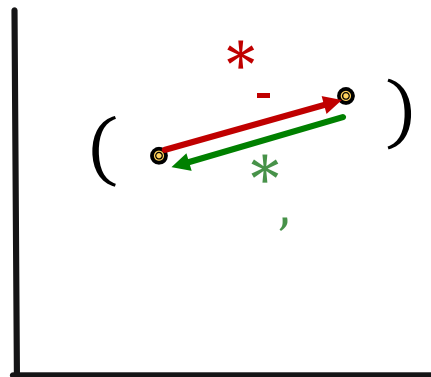
- **Inverse Relations:**

$$R_{\#}(h, \phi) \Rightarrow R_{\$}(\phi, h)$$

§ **Example** : (Advisor, Advisee)

- **TransE** can model inverse relations  $\ddot{u}$

§ ( + \* , = ), we can set \* ) = -\*(



# Composition in TransE

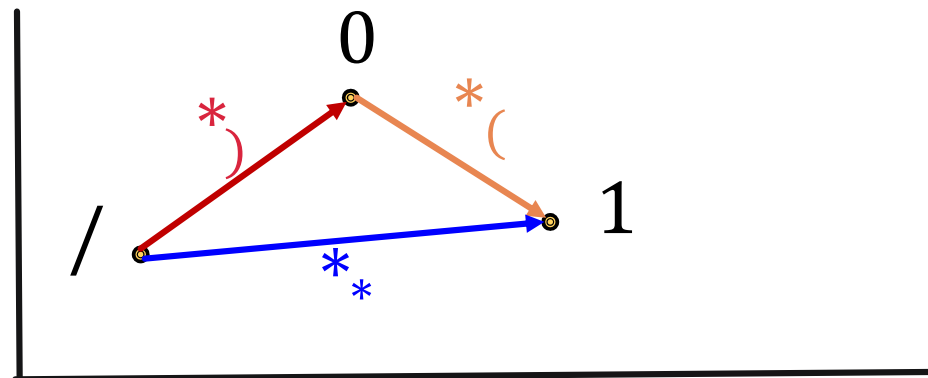
- Composition (Transitive) Relations:

$$r_{\$}(7, 8) \wedge r_{\#}(8, :) \Rightarrow r_{\%}(7, :) \quad \forall 7, 8, :$$

§ Example: My mother's husband is my father.

- TransE can model composition relations

$$r_{\%} = r_{\$} + r_{\#}$$



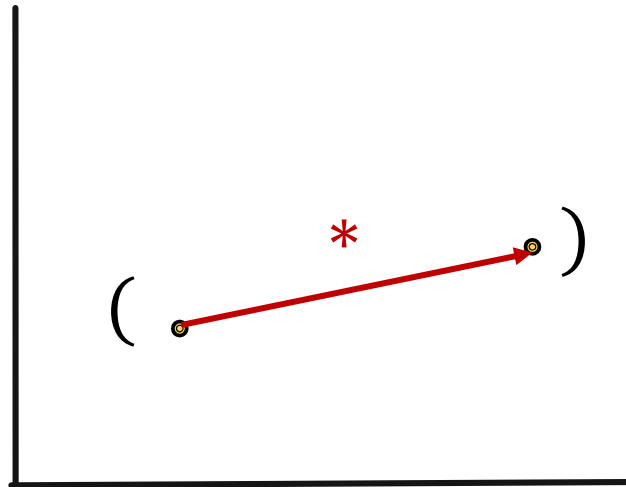
# Limitations of TransE: Symmetric Relations

## i Symmetric Relations:

$$r(h, t) \Rightarrow r(t, h) \quad \forall h, t$$

§ Example: Family, Roommate

## i TransE cannot model symmetric relations $\hat{U}$ only if $(\mathbf{r} = \mathbf{0}, \mathbf{h} = \mathbf{t})$



For all  $h, t$  that satisfy  $r(h, t)$ ,  $r(t, h)$  is also True, which means  $\|\mathbf{h} + \mathbf{r} - \mathbf{t}\| = 0$  and  $\|\mathbf{t} + \mathbf{r} - \mathbf{h}\| = 0$ . Then  $\mathbf{r} = \mathbf{0}$  and  $\mathbf{h} = \mathbf{t}$ , however  $h$  and  $t$  are two different entities and should be mapped to different locations.



# Limitations of TransE: 1-to-N Relations

## i 1-to-N Relations:

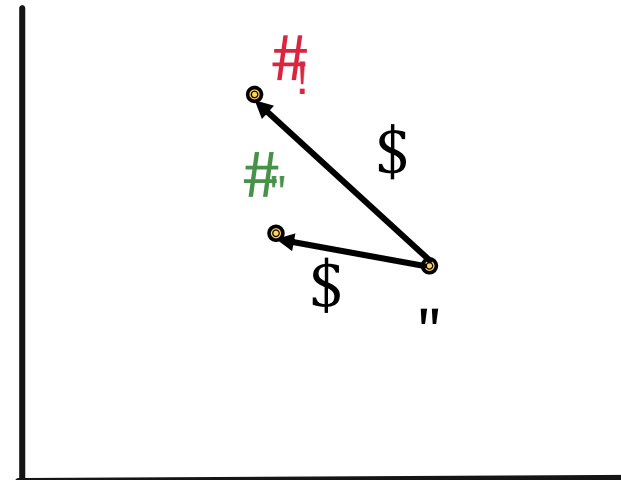
§ **Example:**  $(h, r, t_j)$  and  $(h, r, t_c)$  both exist in the knowledge graph, e.g.,  $r$  is “StudentsOf”

## i TransE cannot model 1-to-N relations $\hat{u}$

§  $t_j$  and  $t_c$  will map to the same vector, although they are different entities

i  $t_j = h + r = t_c$

i  $t_j \neq t_c$  **contradictory!**



# KG Completion Methods

Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^l$	$\hat{u}$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$	$\hat{u}$
TransR	$-\ M_r \mathbf{h} + \mathbf{r} - M_r \mathbf{t}\ $	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^l,$ $\mathbf{r} \in \mathbb{R}^\#,$ $M_r \in \mathbb{R}^{\# \times l}$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$
DistMult	$\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^l$	$\ddot{u}$	$\hat{u}$	$\hat{u}$	$\hat{u}$	$\ddot{u}$
Complex	$\text{Re}(\langle \mathbf{h}, \mathbf{r}, \bar{\mathbf{t}} \rangle)$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{C}^l$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$	$\hat{u}$	$\ddot{u}$
RotateE	$-\ \mathbf{h} \circ \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{C}^l$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$	$\ddot{u}$

**Questions?**