

DSC250: Advanced Data Mining

Knowledge Graphs (KGs)

Zhiting Hu

Lecture 14, Feb 20, 2025

UC San Diego

HALICIOĞLU DATA SCIENCE INSTITUTE

Outline

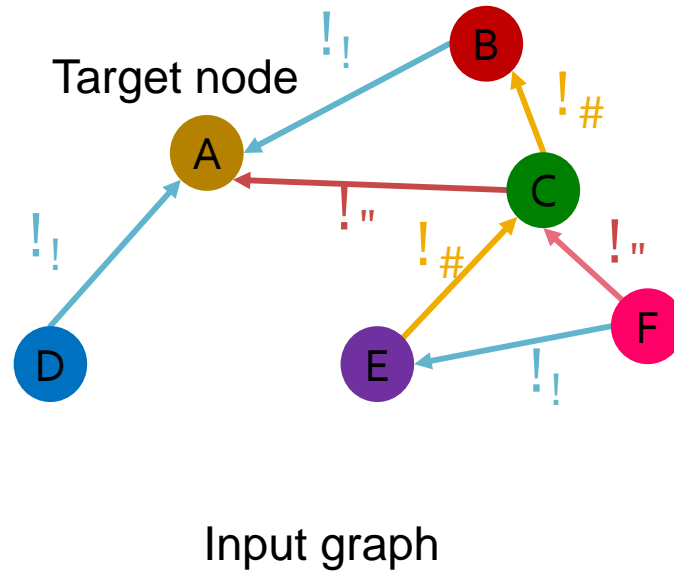
- Knowledge graphs
- Presentation
 - Lingyi Yang, Jiayue Yuan: "Generative Agents: Interactive Simulacra of Human Behavior"
 - Yuyuan Wu, Zijie Feng: "Data-efficient Fine-tuning for LLM-based Recommendation"
 - Ivy Nguyen, Joshua Chuang: "Prediction of COVID-19 cases by multifactor driven long short-term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

Outline

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

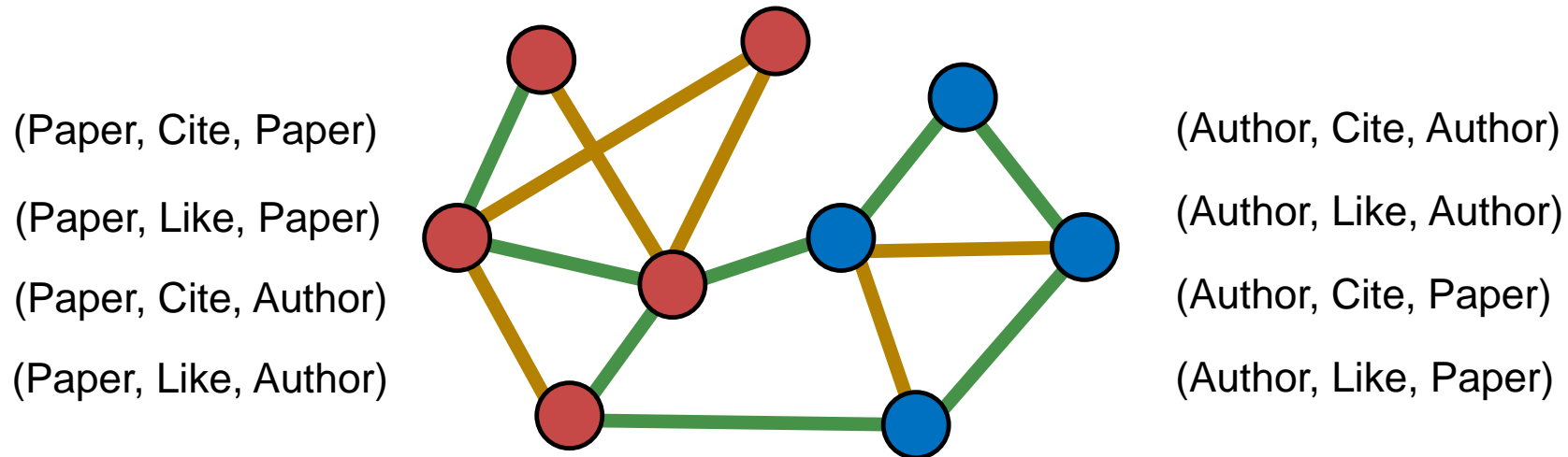
Heterogeneous Graphs

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



Heterogeneous Graphs

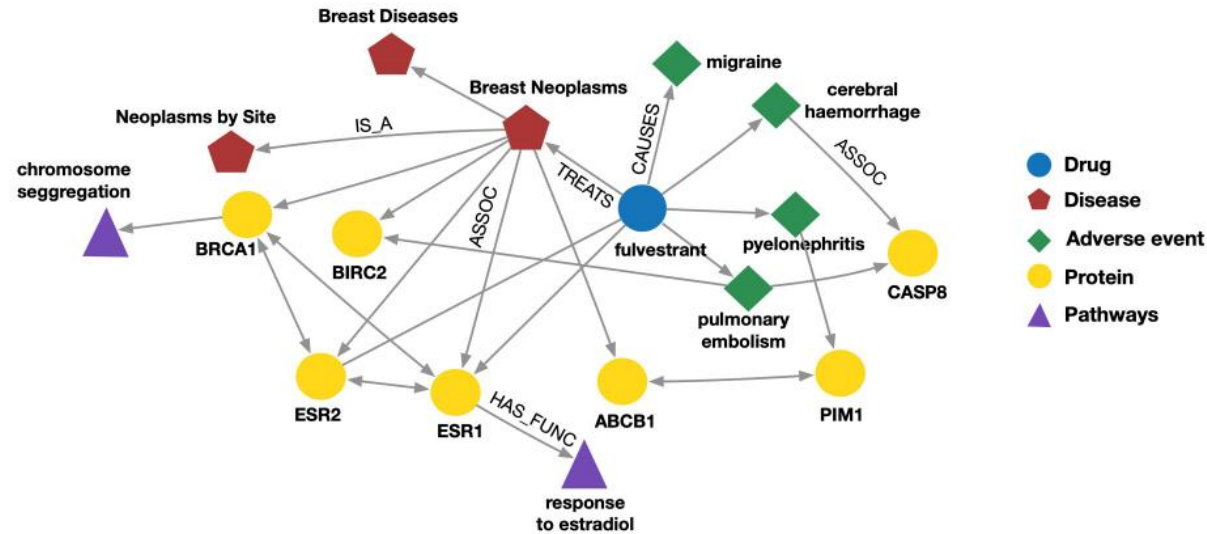
8 possible relation types!



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



Biomedical Knowledge Graphs

Example node: Migraine

Example relation: (fulvestrant, Treats, Breast Neoplasms)

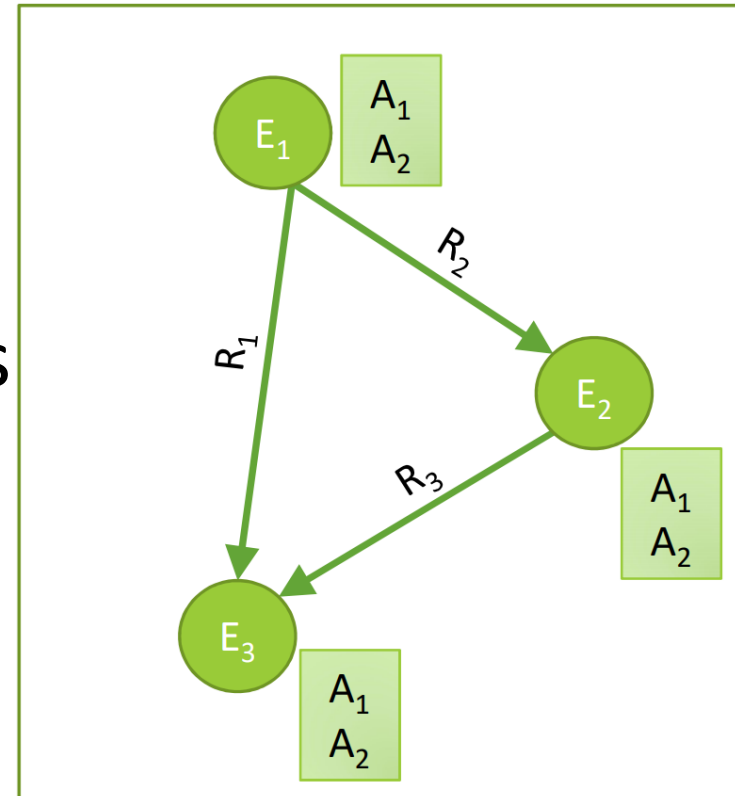
Example node type: Protein

Example edge type: Causes

Knowledge Graph

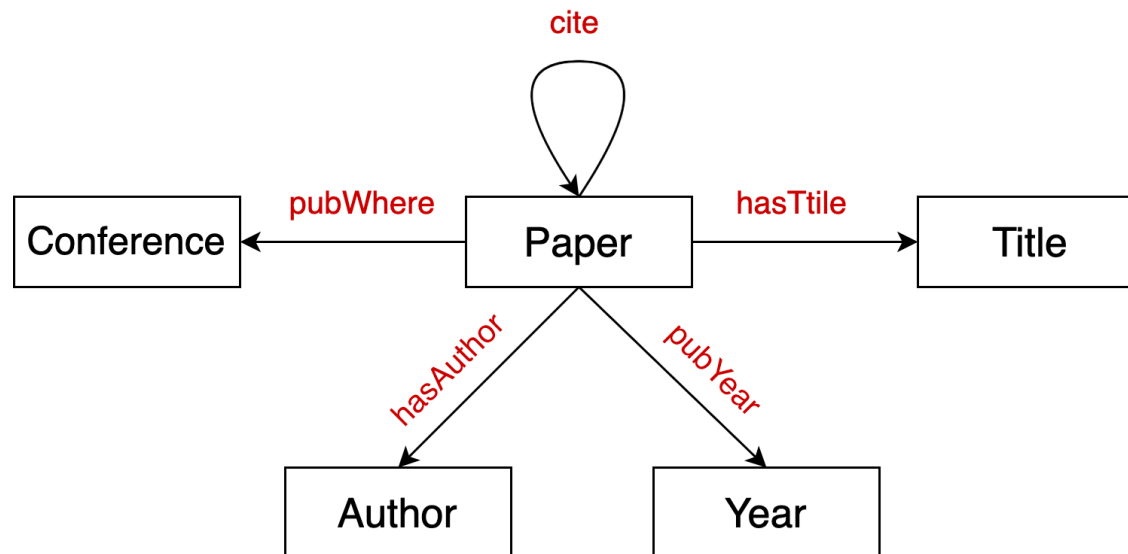
Knowledge in graph form:

- § 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**



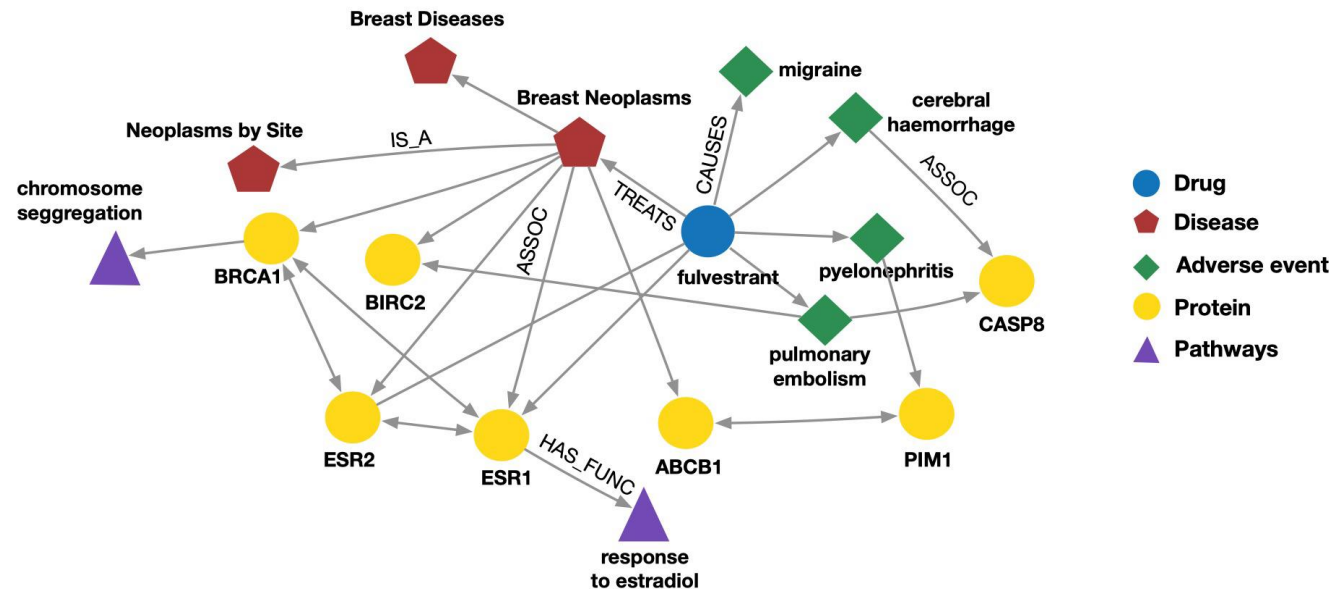
Example: Bibliographic Networks

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



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:

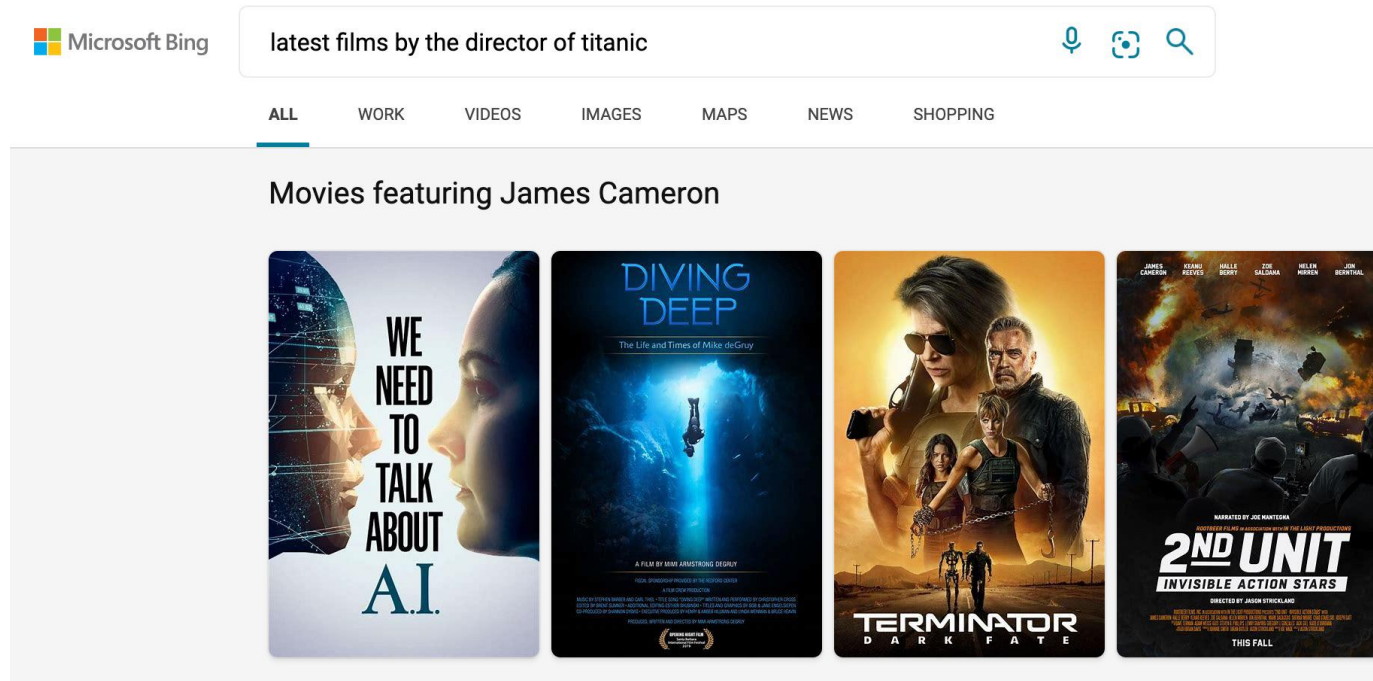


Image credit: Bing

Applications of KGs

i Question answering and conversation agents

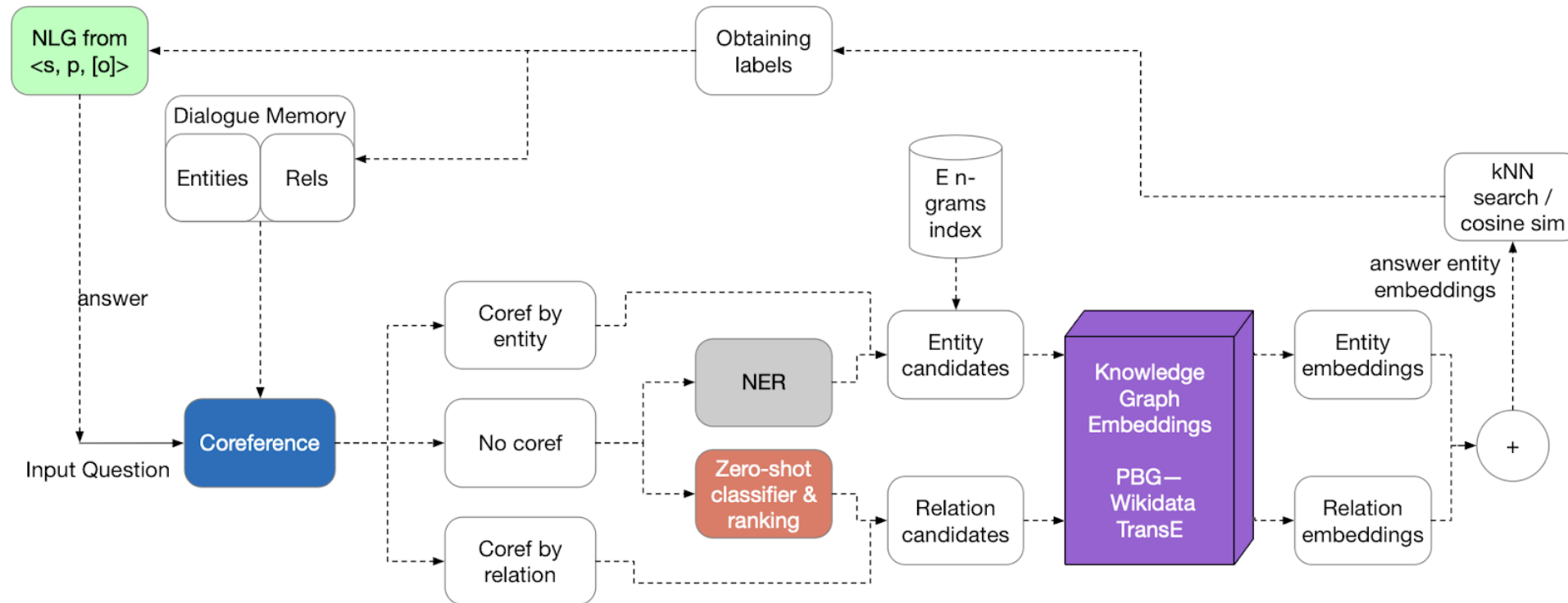


Image credit: [Medium](#)

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

KG Datasets

- | **Publicly available KGs:**

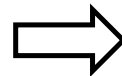
- § FreeBase, Wikidata, Dbpedia, YAGO, NELL, etc.

- | **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

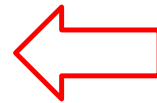


§ ~80 million **entities**

§ ~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

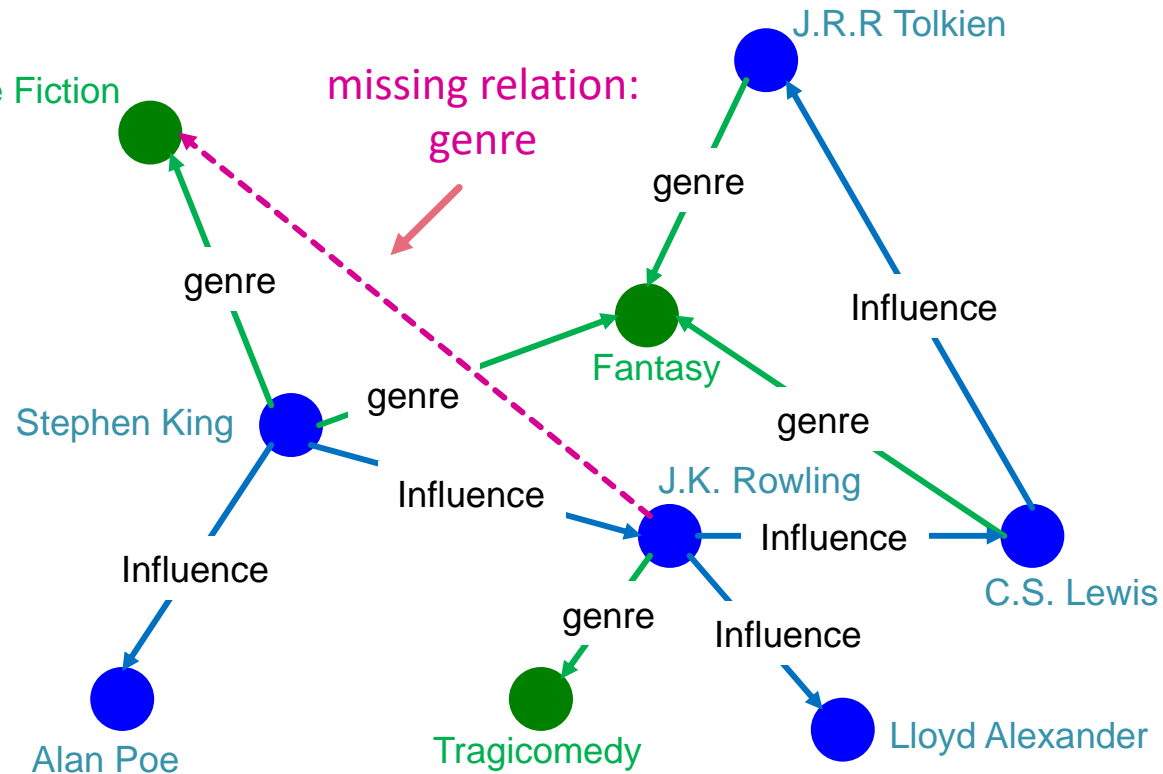
KG Completion Task

Given an enormous KG, can we complete the KG?

§ For a given (**head**, **relation**), we predict missing **tails**.

§ (Note this is slightly different from link prediction task)

Example task: predict the tail “Science Fiction” for (“J.K. Rowling”, “genre”)



KG Representation

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

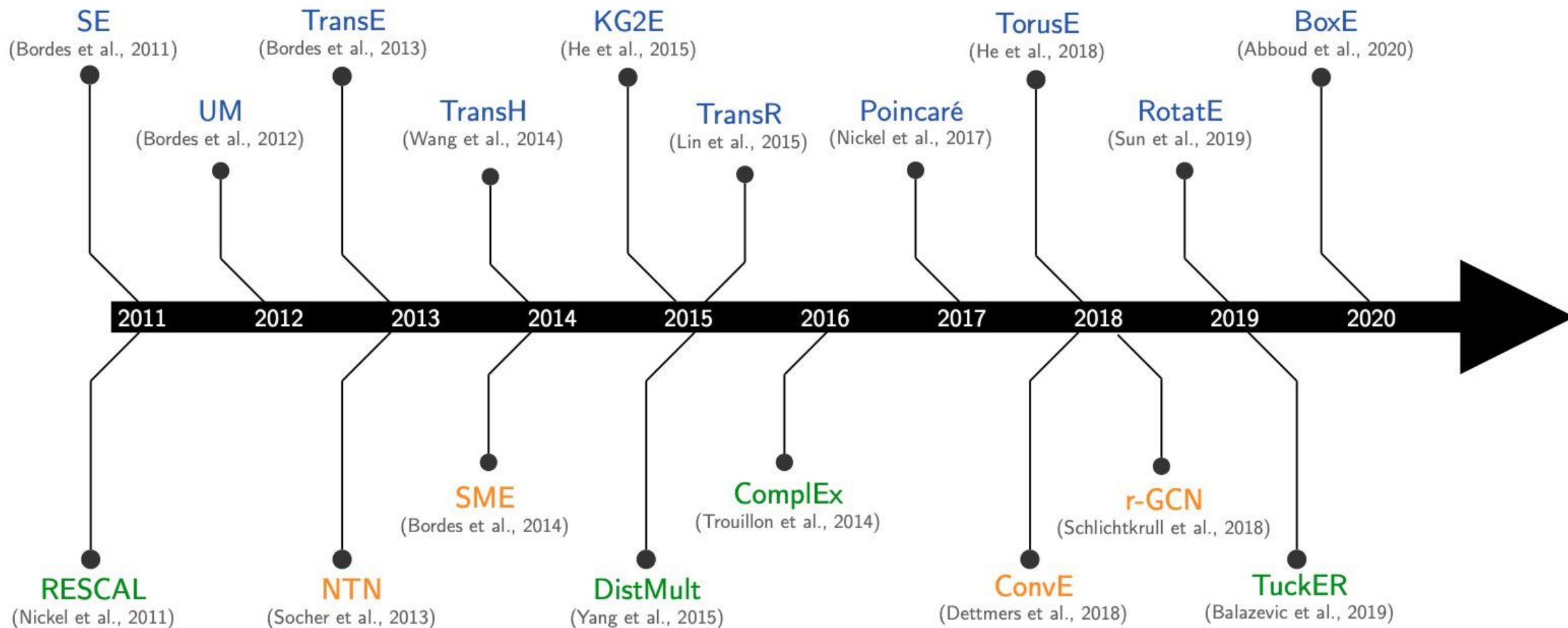
KG Representation

- 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, 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!**
 - § Given a triple (h, r, t) , the goal is that the **embedding of (h, r) should be close** to the **embedding of t**
 - § How to embed (h, r) ?
 - § How to define score $f_r(h, t)$?
 - § Score f_r is high if (h, r, t) exists, else f_r is low

Many KG Embedding Methods



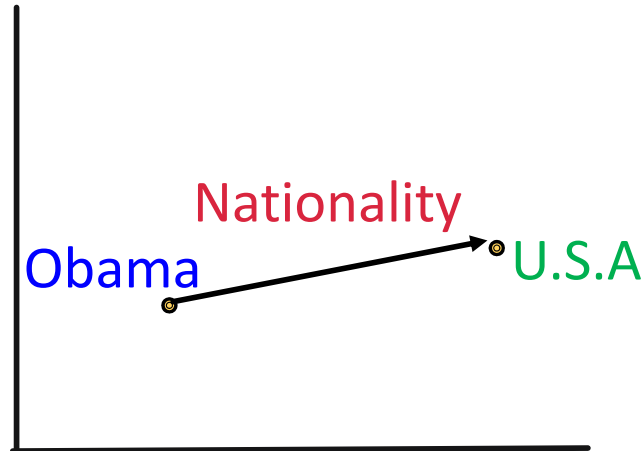
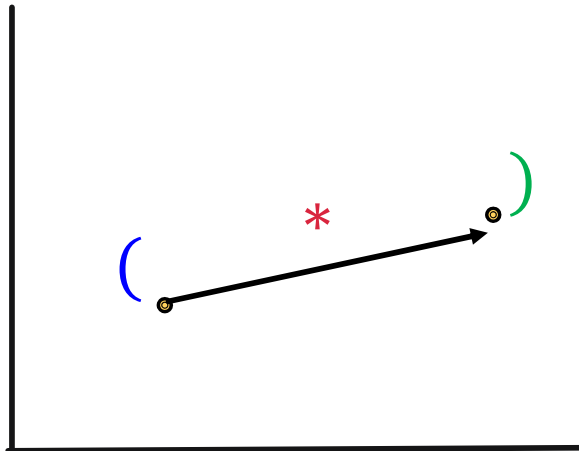
TransE for KG Completion

- 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

- 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

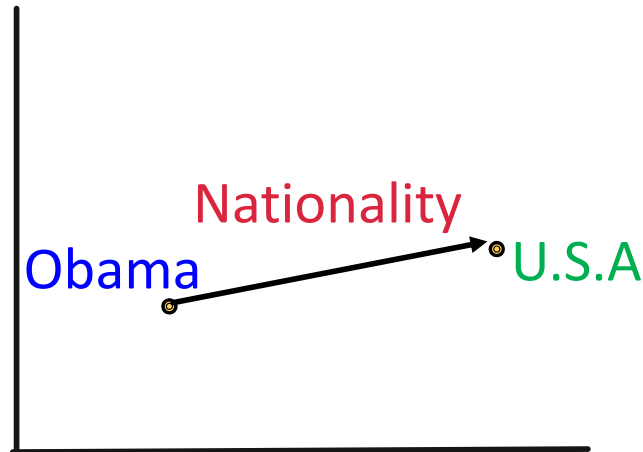
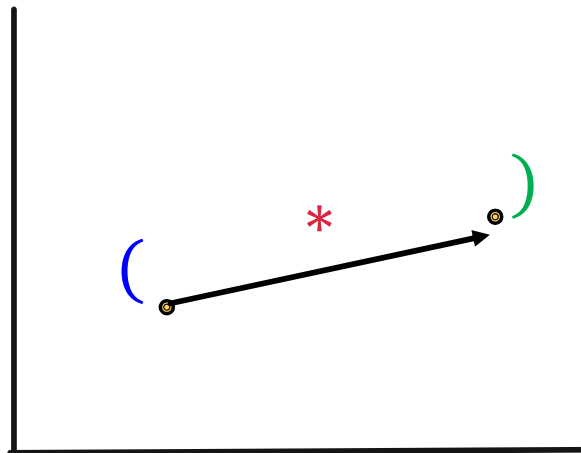
i Intuition: Translation

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



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:**

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

§ **Example:**

§ Symmetric: Family, Roommate

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

i **Inverse Relations:**

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

§ **Example** : (Advisor, Advisee)

i **Composition (Transitive) Relations:**

$$r_r(x, y) \wedge r_l(y, z) \Rightarrow r_*(x, z) \quad \forall x, y, z$$

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

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

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

§ **Example**: # is "StudentsOf"

Antisymmetric Relations in TransE

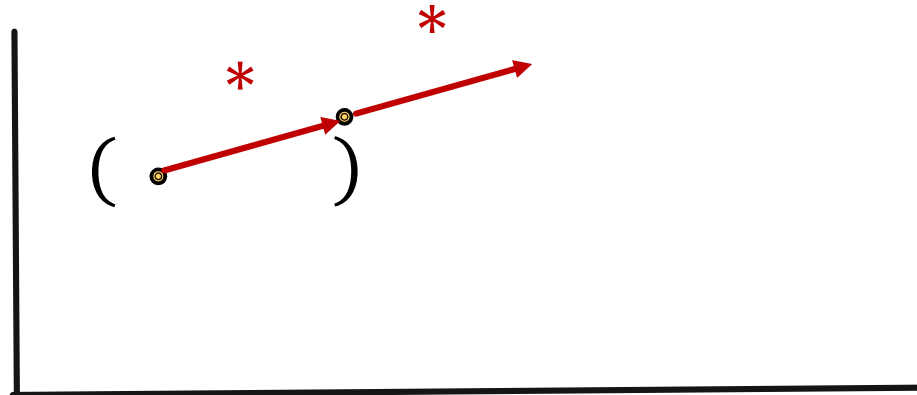
i **Antisymmetric Relations:**

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

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

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

§ (+ * =), but) + * ≠ (



Inverse Relations in TransE

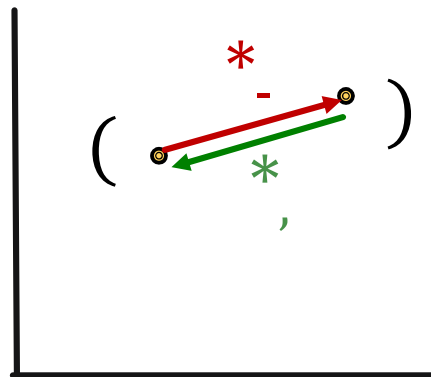
- Inverse Relations:**

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

§ **Example** : (Advisor, Advisee)

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

§ (+ * , =), we can set *) = -*(



Composition in TransE

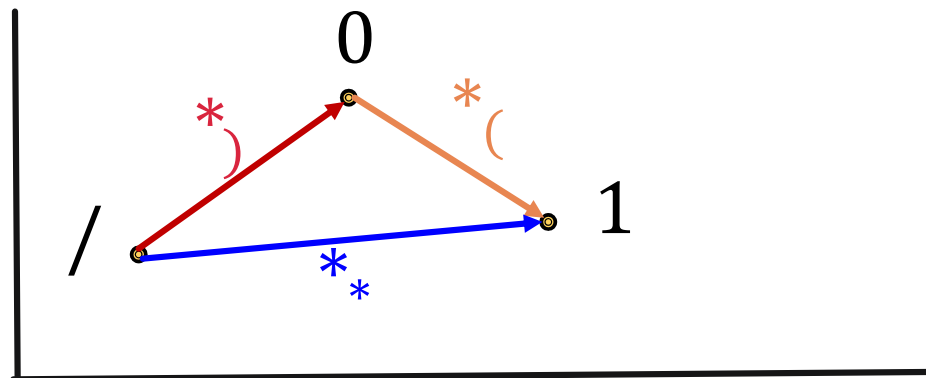
- Composition (Transitive) Relations:

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

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

- TransE can model composition relations

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



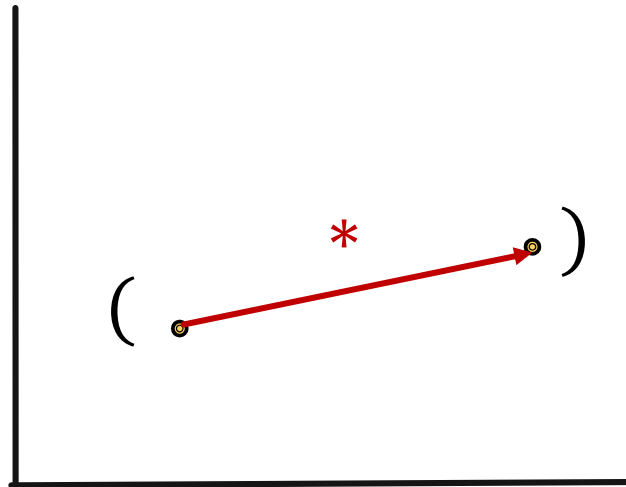
Limitations of TransE: Symmetric Relations

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 $(\mathbf{h} = \mathbf{0}, \mathbf{t} = \mathbf{0})$



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:

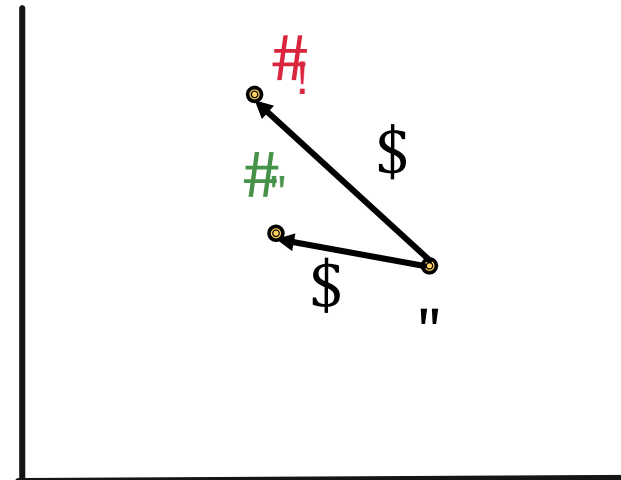
§ **Example:** (h, r, t_j) and (h, r, t_c) both exist in the knowledge graph, e.g., r is “StudentsOf”

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_{\mathbf{r}} \mathbf{h} + \mathbf{r} - M_{\mathbf{t}} \mathbf{t}\ $	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^l,$ $\mathbf{r} \in \mathbb{R}^{\#},$ $M_{\mathbf{r}}, M_{\mathbf{t}} \in \mathbb{R}^{\# \times l}$	\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}

Outline

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

Reasoning over KGs

- i **Goal:**

- § How to perform multi-hop reasoning over KGs?

- i **Reasoning over Knowledge Graphs**

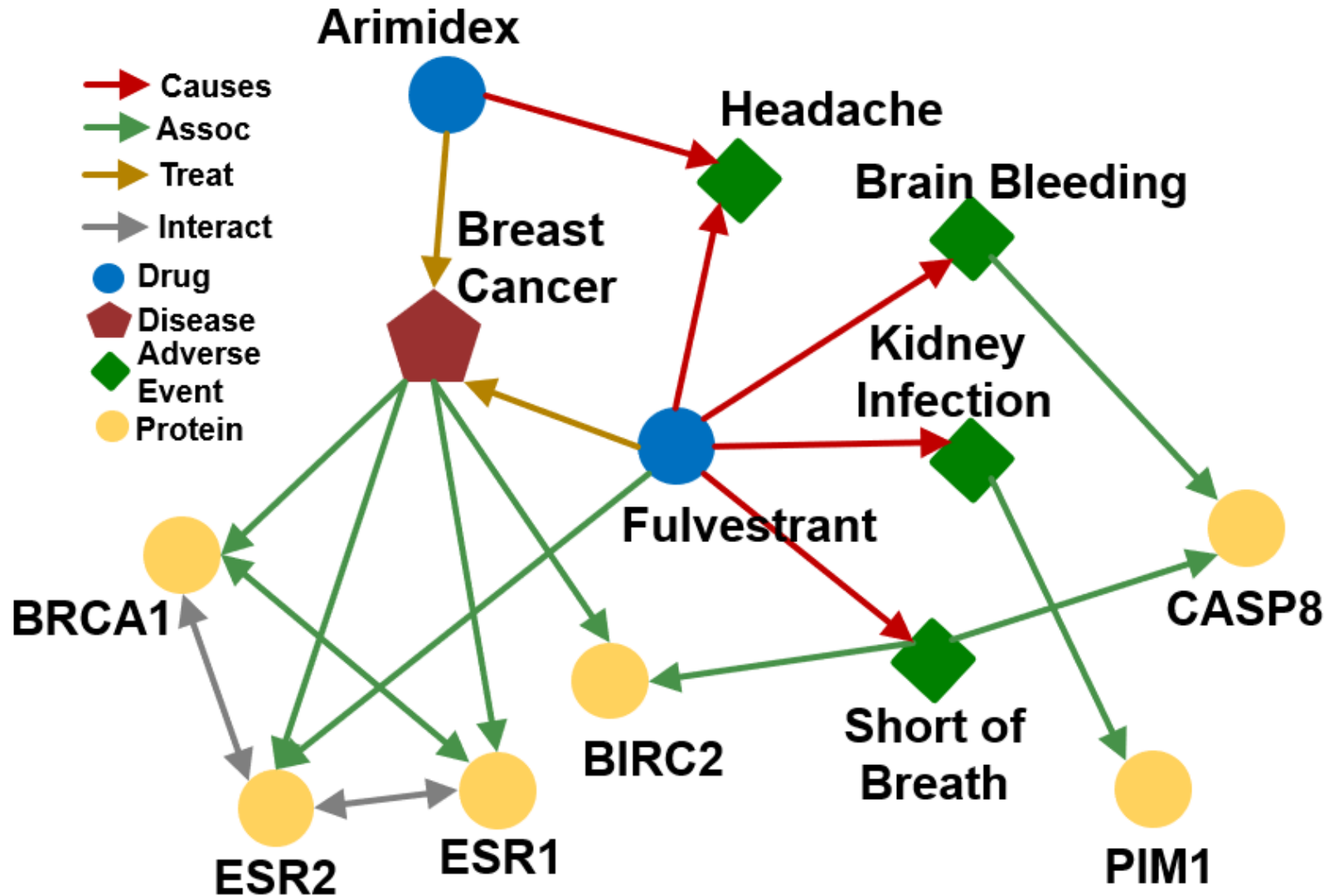
- § Answering multi-hop queries

- § Path Queries

- § Conjunctive Queries

- § Query2Box

Example KG: Biomedicine



Predictive Queries on KG

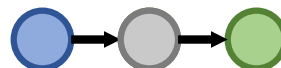
Can we do multi-hop reasoning, i.e., **answer complex queries on an incomplete, massive KG?**

Query Types	Examples: Natural Language Question, Query
One-hop Queries	What adverse event is caused by Fulvestrant? (e:Fulvestrant, (r:Causes))
Path Queries	What protein is associated with the adverse event caused by Fulvestrant? (e:Fulvestrant, (r:Causes, r:Assoc))
Conjunctive Queries	What is the drug that treats breast cancer and caused headache? ((e:BreastCancer, (r:TreatedBy)), (e:Migraine, (r:CausedBy)))

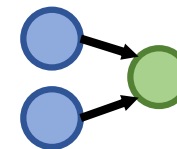
In this lecture, we only focus on answering **queries** on a KG!
The notation will be detailed next.



One-hop Queries



Path Queries

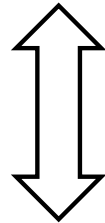


Conjunctive Queries

Predictive One-hop Queries

i We can formulate knowledge graph completion problems as answering one-hop queries.

i **KG completion:** Is link (h, r, o) in the KG?



i **One-hop query:** Is o an answer to query (h, r)?

§ **For example:** What side effects are caused by drug Fulvestrant?

Path Queries

- Generalize one-hop queries to path queries by **adding more relations on the path**.

- An n -hop path query q can be represented by

$$q = (v_1, (r_1, \dots, r_n))$$

§ v_1 is an “anchor” entity,

§ Let answers to q in graph G be denoted by $[[q]]_G$.

Path Queries

- Generalize one-hop queries to path queries by **adding more relations on the path**.

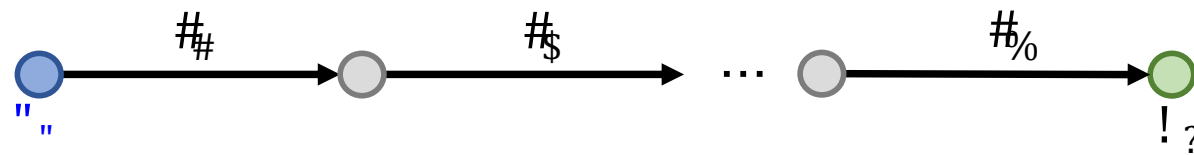
- An n -hop path query q can be represented by

$$q = (v_1, (r_1, \dots, r_n))$$

§ v_1 is an “anchor” entity,

§ Let answers to q in graph G be denoted by $[[q]]_G$.

Query Plan of q :

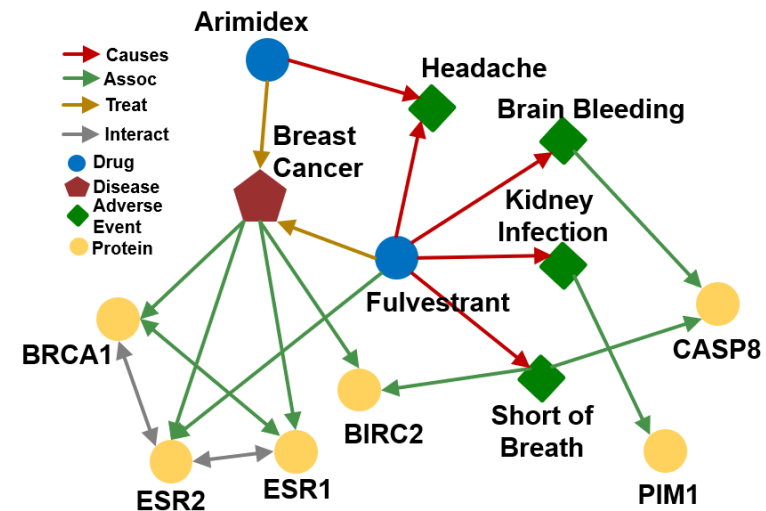
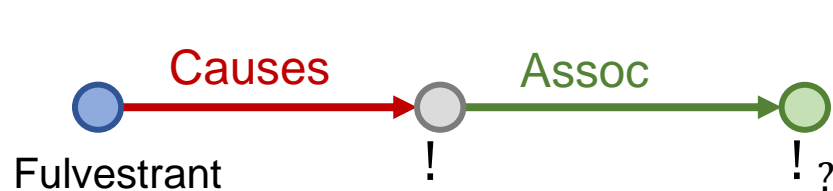


Query plan of path queries is a chain.

Path Queries

Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”

- ! is e:Fulvestrant
- (\$_1, \$_2) is (r:Causes, r:Assoc)
- Query: (e:Fulvestrant, (r:Causes, r:Assoc))

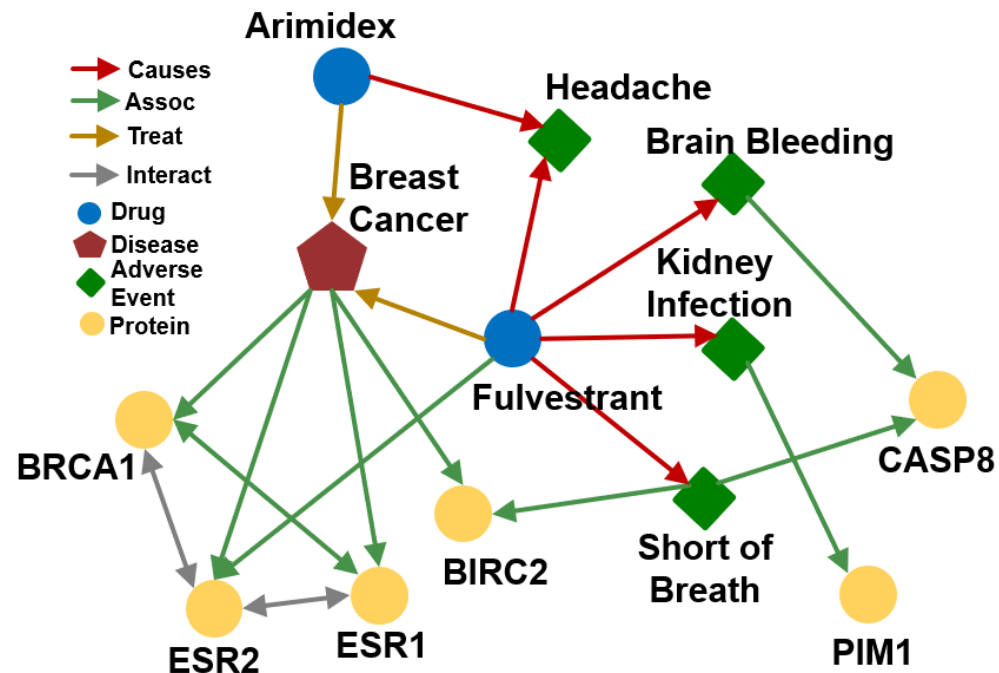


Path Queries

Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”

Query: (e:Fulvestrant, (r:Causes, r:Assoc))

Given a KG, how to answer a path query?



Traversing Knowledge Graphs

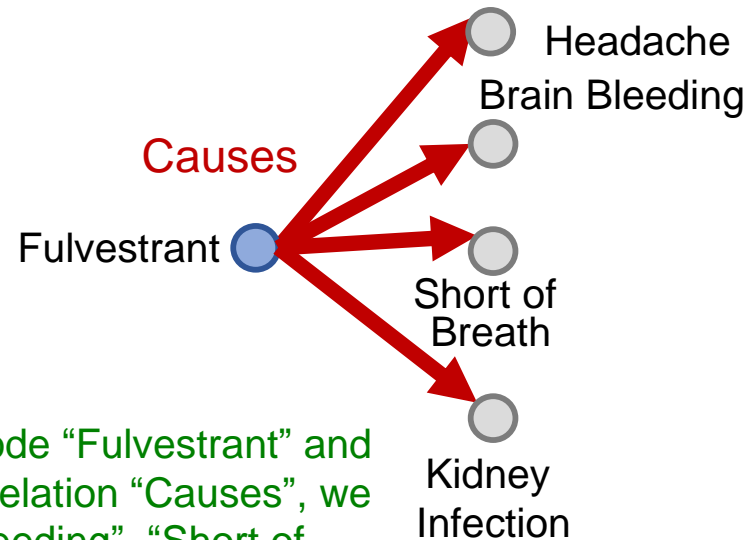
- i We answer path queries by traversing the KG:
“What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- i Query: (e:Fulvestrant, (r:Causes, r:Assoc))

Fulvestrant 

Start from the
anchor node
(Fulvestrant).

Traversing Knowledge Graphs

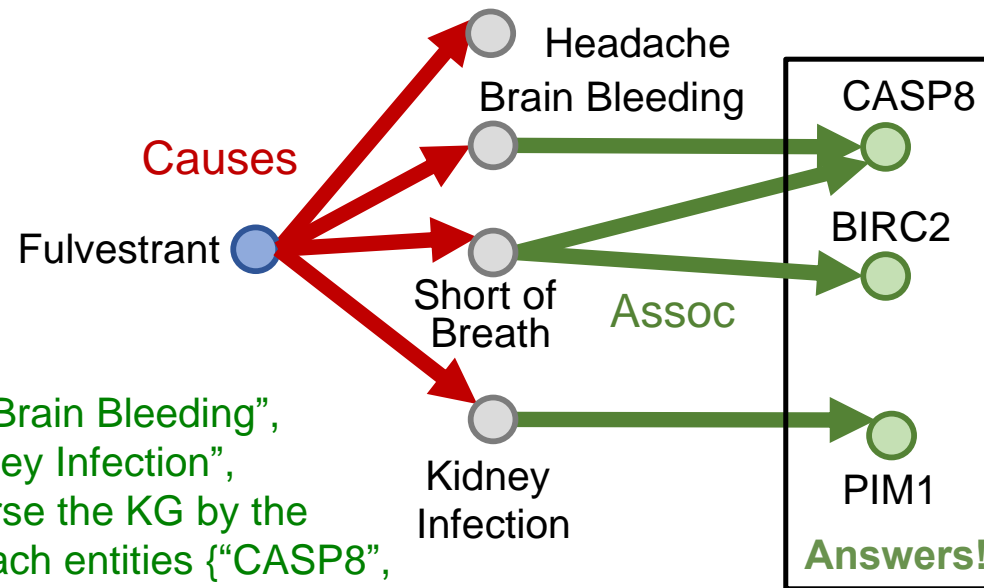
- ! We answer path queries by traversing the KG:
“What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- ! Query: (e:Fulvestrant, (r:Causes, r:Assoc))



Start from the anchor node “Fulvestrant” and traverse the KG by the relation “Causes”, we reach entities {“Brain Bleeding”, “Short of Breath”, “Kidney Infection”, “Headache”}.

Traversing Knowledge Graphs

- ! We answer path queries by traversing the KG:
“What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- ! Query: (e:Fulvestrant, (r:Causes, r:Assoc))



Start from the nodes {“Brain Bleeding”, “Short of Breath”, “Kidney Infection”, “Headache”} and traverse the KG by the relation “Assoc”, we reach entities {“CASP8”, “BIRC2”, “PIM1”}. These are the answers.

However, KGs are incomplete

- | **Answering queries seems easy: Just traverse the graph.**
- | **But KGs are incomplete and unknown:**
 - § Many relations between entities are missing or are incomplete
 - § For example, we lack all the biomedical knowledge
 - § Enumerating all the facts takes non-trivial time and cost, we cannot hope that KGs will ever be fully complete
- | **Due to KG incompleteness, one is not able to identify all the answer entities**

Can KG Completion Help?

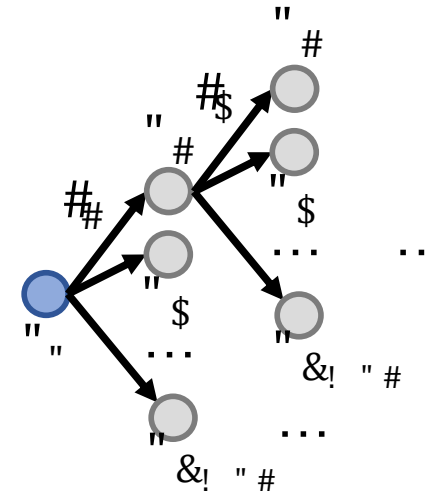
Can we first do KG completion and then traverse the completed (probabilistic) KG?

Can KG Completion Help?

Can we first do KG completion and then traverse the completed (probabilistic) KG?

- ! **No!** The “completed” KG is a **dense graph!**
 - § Most (h, r, t) triples (edge on KG) will have some non-zero probability.

- ! Time complexity of traversing a dense KG is exponential as a function of the path length, $\sim (|R| \cdot |I|)^L$



Task: Predictive Queries

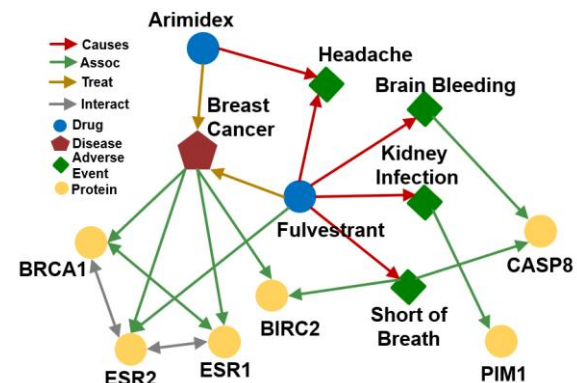
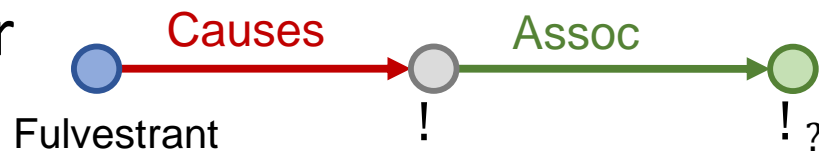
- i We need a way to answer path-based queries over an incomplete knowledge graph.
- i We want our approach to implicitly impute and account for the incomplete KG.

Task: Predictive Queries

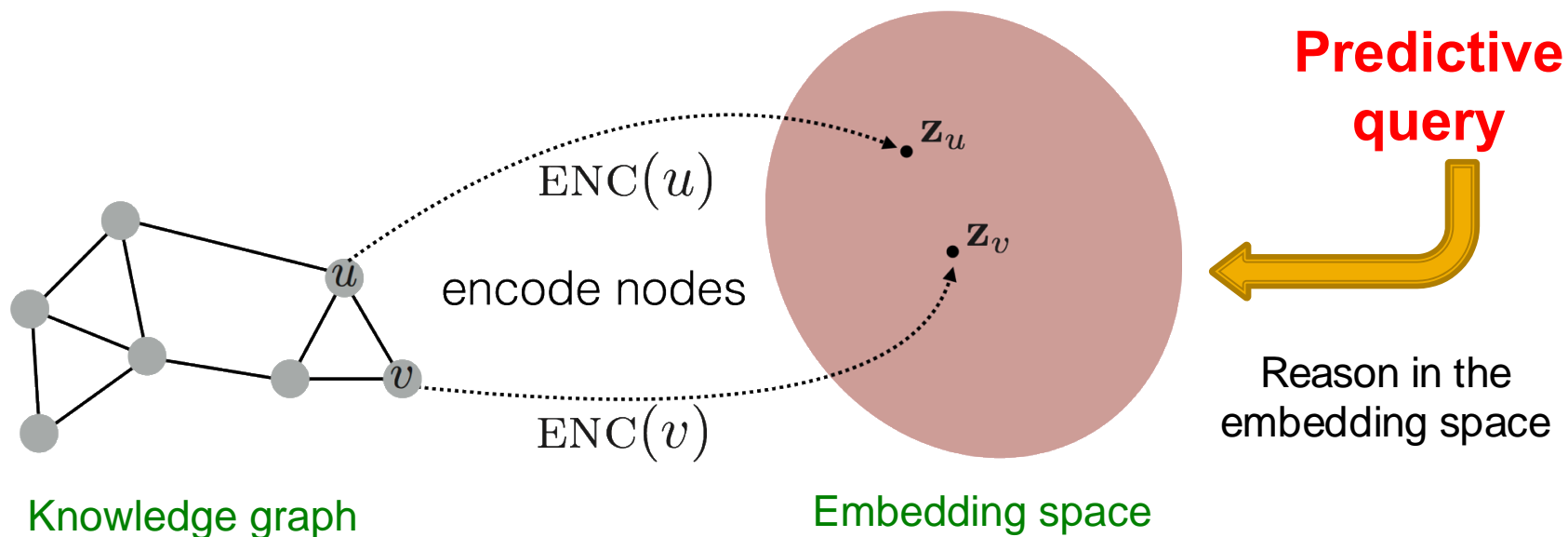
- i We need a way to answer path-based queries over an incomplete knowledge graph.
- i We want our approach to implicitly impute and account for the incomplete KG.
- i **Task: Predictive queries**

§ Want to be able to answer arbitrary queries while implicitly imputing for the missing information

§ **Generalization of the link prediction task**



A General Idea



Map queries into embedding space. **Learn to reason in that space**

- Embed query into a single **point** in the Euclidean space: answer nodes are close to the query.
- Query2Box**: Embed query into a hyper-rectangle (**box**) in the Euclidean space: answer nodes are enclosed in the box.

[[Embedding Logical Queries on Knowledge Graphs](#). Hamilton, et al., NeurIPS 2018]

[[Query2box: Reasoning over Knowledge Graphs in Vector Space Using Box Embeddings](#). Ren, et al., ICLR 2020]

Traversing KG in Vector Space

- ! **Key idea: Embed queries!**

- § Generalize **TransE** to multi-hop reasoning.

Traversing KG in Vector Space

! Key idea: Embed queries!

§ Generalize **TransE** to multi-hop reasoning.

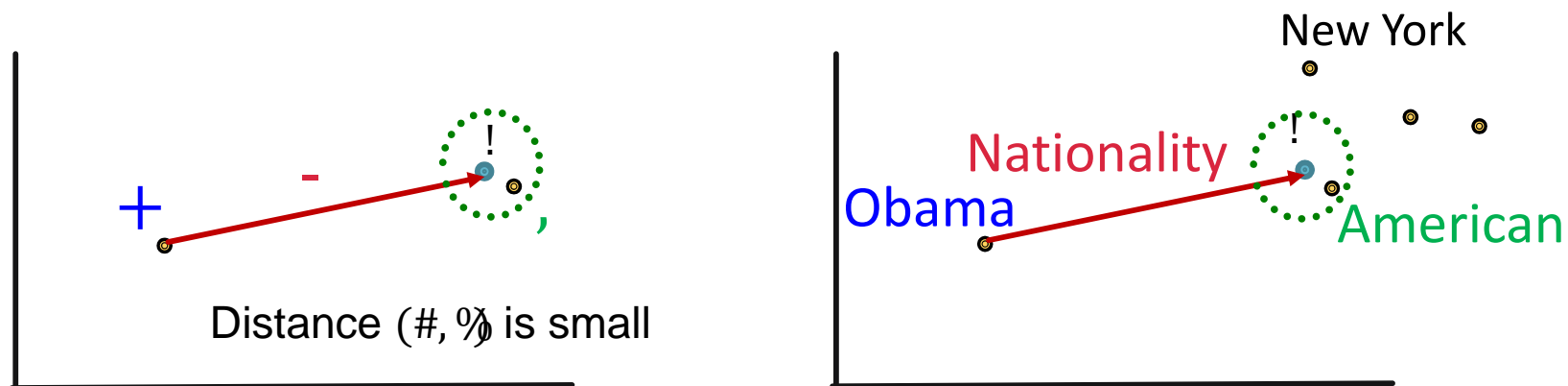
§ **Recap: TransE:** Translate ! to " using # with score function $f(h, r) = -\|h + r - t\|$.

§ Another way to interpret this is that:

§ **Query embedding:** $q = h + r$

§ Goal: **query embedding** q is **close** to the **answer embedding** t

$$f_i(t) = -\|q - t\|$$

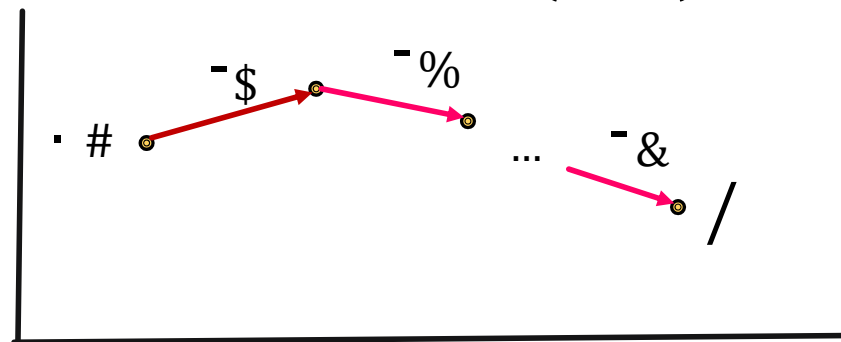


Traversing KG in Vector Space

! Key idea: Embed queries!

§ Generalize TransE to multi-hop reasoning.

Given a path query $q = (0, (1), \dots, 1^*)$,



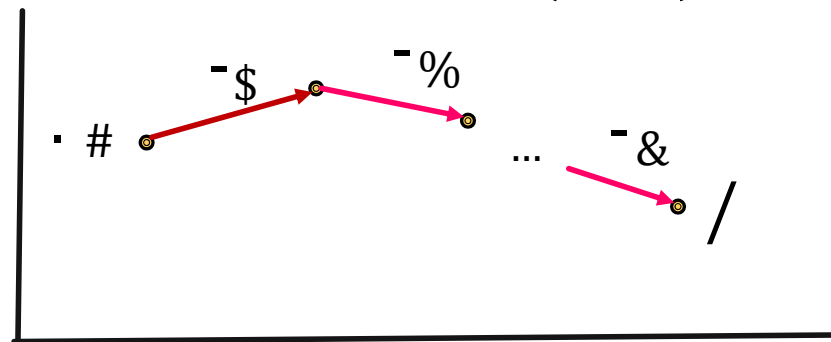
$$2 = 3_1 + 5_1 + \dots + 5_{\#}$$

Traversing KG in Vector Space

i Key idea: Embed queries!

§ Generalize TransE to multi-hop reasoning.

Given a path query $q = (0, (1), \dots, 1^*)$,



$$2 = 3_1 + 5_1 + \dots + 5_{\#}$$

- i The embedding process **only involves vector addition**, **independent of # entities** in the KG!

Traversing KG in Vector Space

Embed path queries in vector space.

- | **Question:** “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- | **Query:** (e:Fulvestrant, (r:Causes , r:Assoc))

Follow the query plan:

Query Plan

Embedding Process

Fulvestrant ●

Fulvestrant ●

Traversing KG in Vector Space

Embed path queries in vector space.

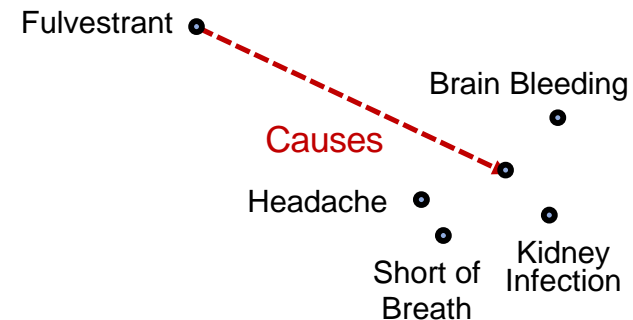
- Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- Query: (e:Fulvestrant, (r:Causes , r:Assoc))

Follow the query plan:

Query Plan



Embedding Process



Traversing KG in Vector Space

Embed path queries in vector space.

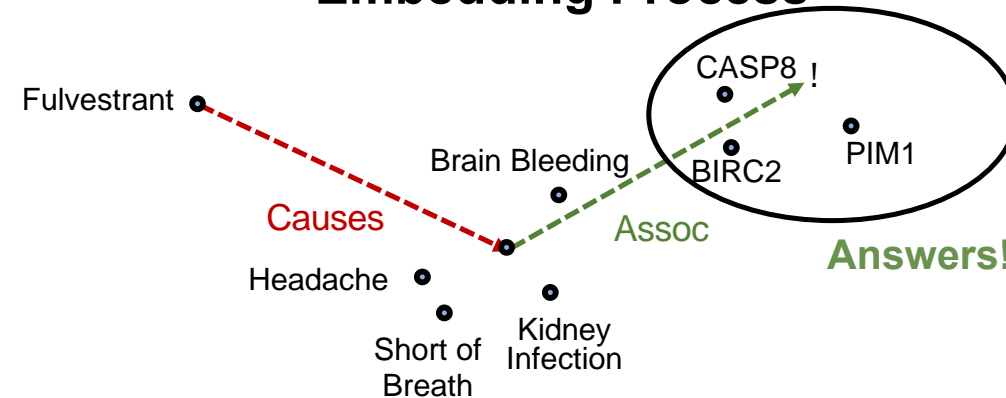
- Question: “What proteins are *associated* with adverse events *caused* by *Fulvestrant*?”
- Query: (e:Fulvestrant, (r:Causes , r:Assoc))

Follow the query plan:

Query Plan



Embedding Process



Traversing KG in Vector Space

Insights:

- i We can train **TransE** to optimize knowledge graph completion objective
- i Since **TransE** can naturally handle **compositional relations**, it can handle path queries by translating in the latent space **for multiple hops using addition of relation embeddings.**

Questions?