



Insight



# Identifying Equivalent Relation Paths in Knowledge Graphs

Sameh K. Mohamed, Emir Muñoz, Vít Nováček, Pierre-Yves Vandenbussche

Language, Data and Knowledge (LDK) Conference. Galway, Ireland. 18~21 June 2017

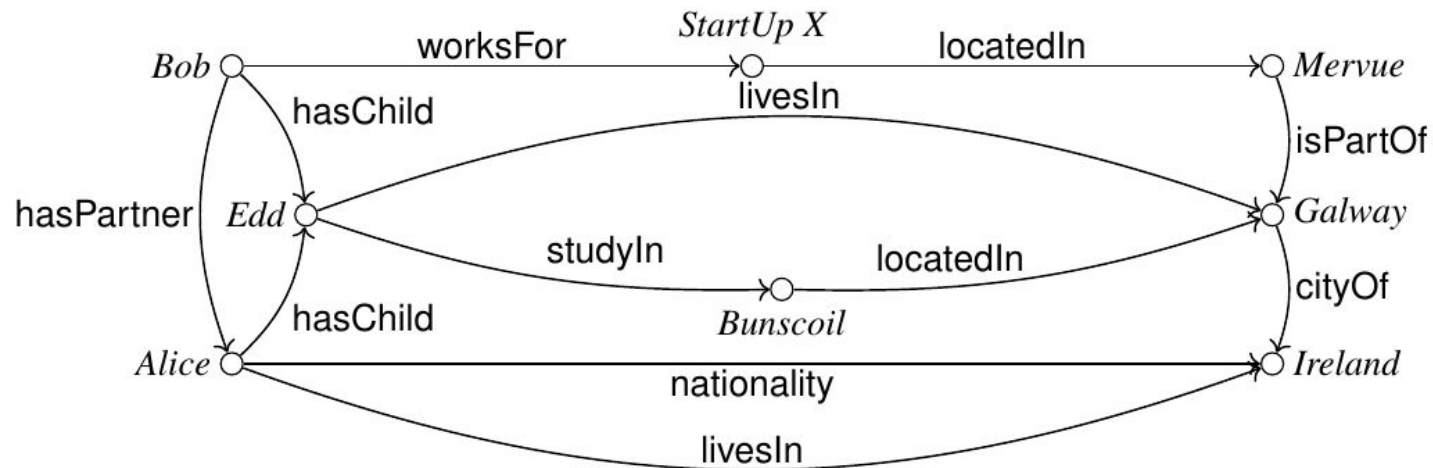


# Agenda

- **Knowledge Graphs**
- **Relation Path Equivalence**
- **Search of  $\Delta$ -Equivalences**
  - Path extension Extraction
  - Subgraph extraction
  - Building connecting Paths
  - Ranking candidate paths
- **Experiments & Results**
- **Conclusions**
- **Future work**

# Knowledge Graphs (KGs)

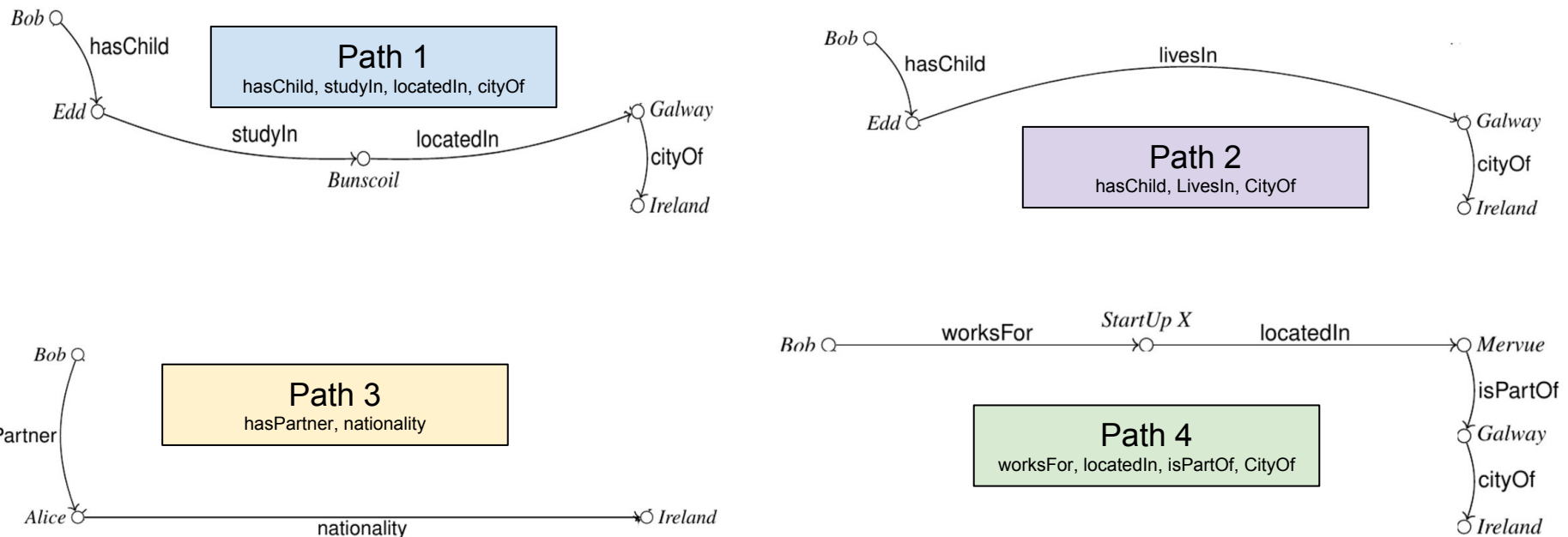
- Graph representation of knowledge bases, where entities are represented by nodes and relations are represented by edges between them
- Commonly, KGs model facts as Subject-Predicate-Object (SPO) triples.
- A sequence of connected edges in a knowledge graph is called *path*



**Figure 1:** Example of a knowledge graph about people living in Galway, Ireland

# Knowledge Graph Paths

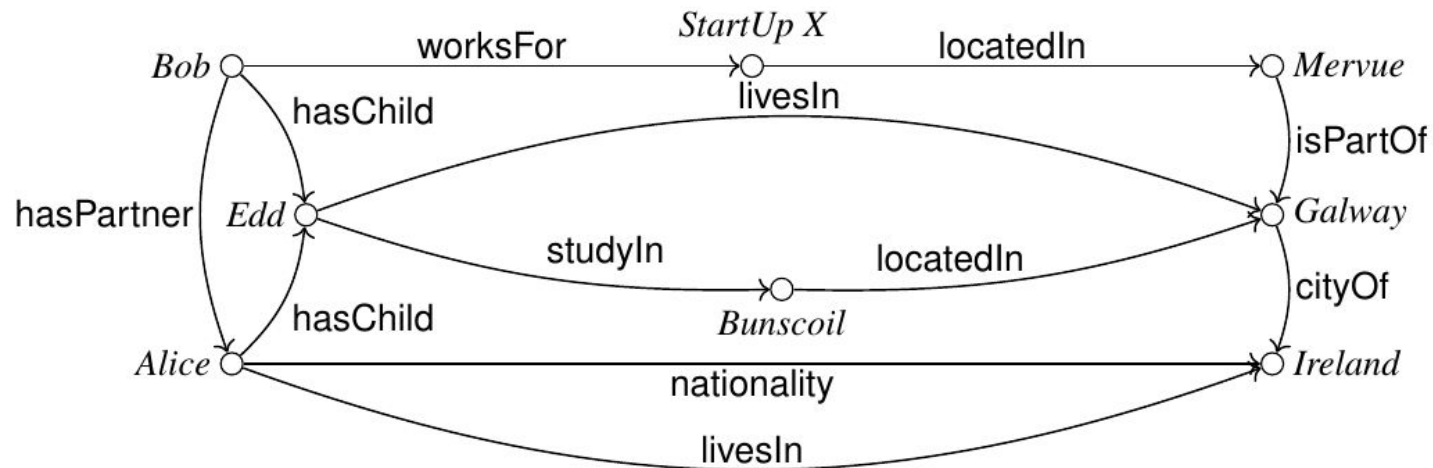
- Relation path is a graph path represented as a sequence of its relations.
- Paths can be navigated in both directions using inverses of relations



**Figure 2:** Sample of paths connecting Bob to Ireland from knowledge graph at Fig. 1

# Knowledge Graph Paths

- **Relation paths are an important feature, that can be used for:**
  - **Expressing properties of entities**
  - **Automatic inference of new facts**
  - **Learning new relation rules**



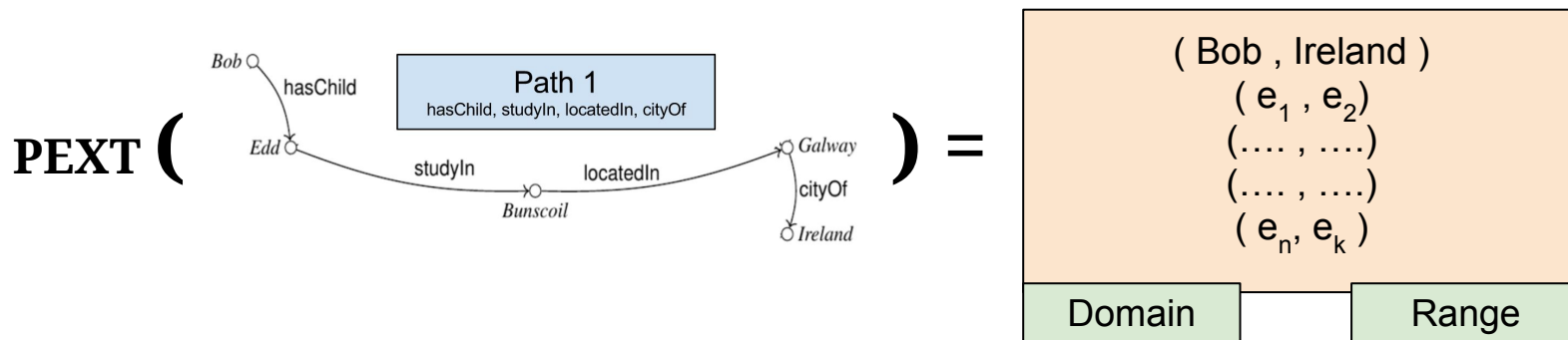
**Figure 1:** Example of a knowledge graph about people living in Galway, Ireland

# Equivalence of Relation Paths

Two paths are equivalent *if and only if* they have the same path extension.

$$P1 \equiv P2 \text{ iff } \mathbf{PEXT}(P1) = \mathbf{PEXT}(P2)$$

where  $\mathbf{PEXT}(P)$  is the set of node pairs that are connected with path P.



Path 1 have the following semantic:

"A person X who has child studying in a school in a city located in a country Y"

# Finding Equivalence of a Relation Path

## SEARCH OF EQUIVALENT RELATION PATHS PROBLEM

**Given:** a knowledge graph  $\mathcal{G}$ , relation path query  $Q$ , integer  $k$ , depth  $d$

**Find:** top- $k$  equivalent relation paths of max. length  $2d$  for  $Q$  in  $\mathcal{G}$  according to a ranking function  $Rank_Q(P)$ , for all  $P \in \mathcal{C}$  set of candidates.

### Challenges:

- Extracting path extension is a complex process
- Finding an equivalent relation path require trying combination of all possible path, which is a complex process.
- Knowledge incompleteness affect representation of paths in knowledge graphs, that equivalent paths can have different extension

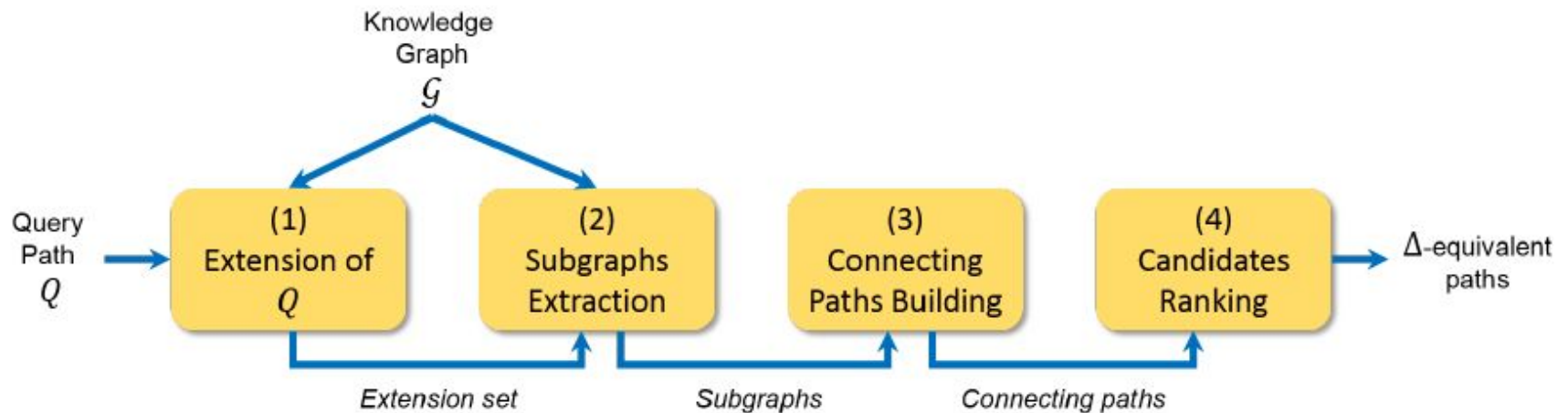
### Proposed Solution:

- A technique for finding approximate equivalences of a relation path using a sample of path extension instances

## Search Approximate Equivalences of a path (1)

We query approximate equivalences of a relation path using a 4 phases procedure:

1. Extracting sample of path extension
2. Extracting Subgraphs of extension instances
3. Build connecting between domain and range nodes of path extension
4. Ranking Candidate connecting paths



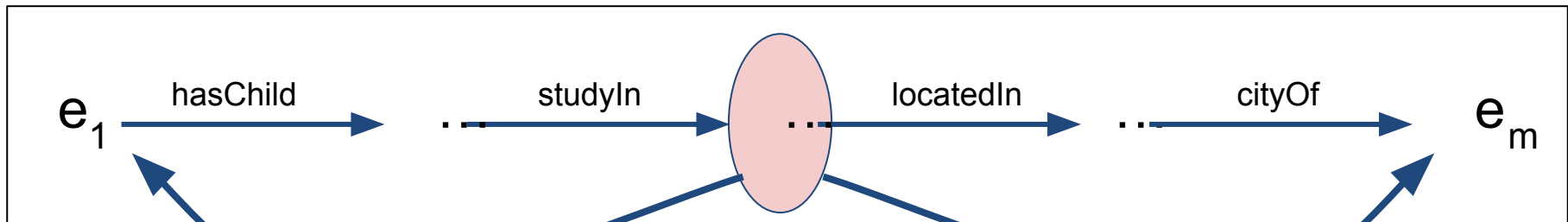
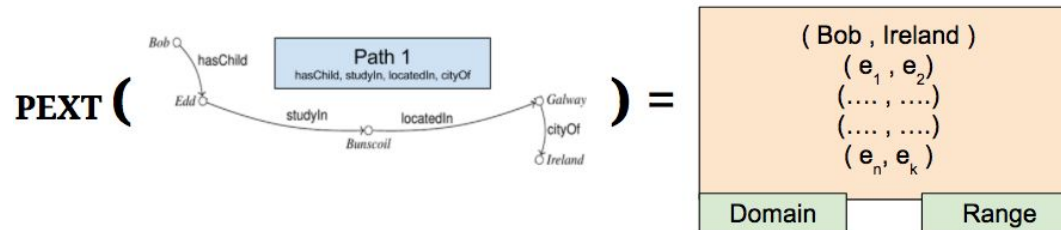
**Figure 2:** Flow diagram of the process of retrieving approximate equivalences of a query path "Q"



## Search Approximate Equivalences of a path (2)

### (1) Extracting sample of path extension

- Finding middle path entities
- Walk from middle points to both domain and range nodes
- Combine reached domain and range nodes starting from same middle node.

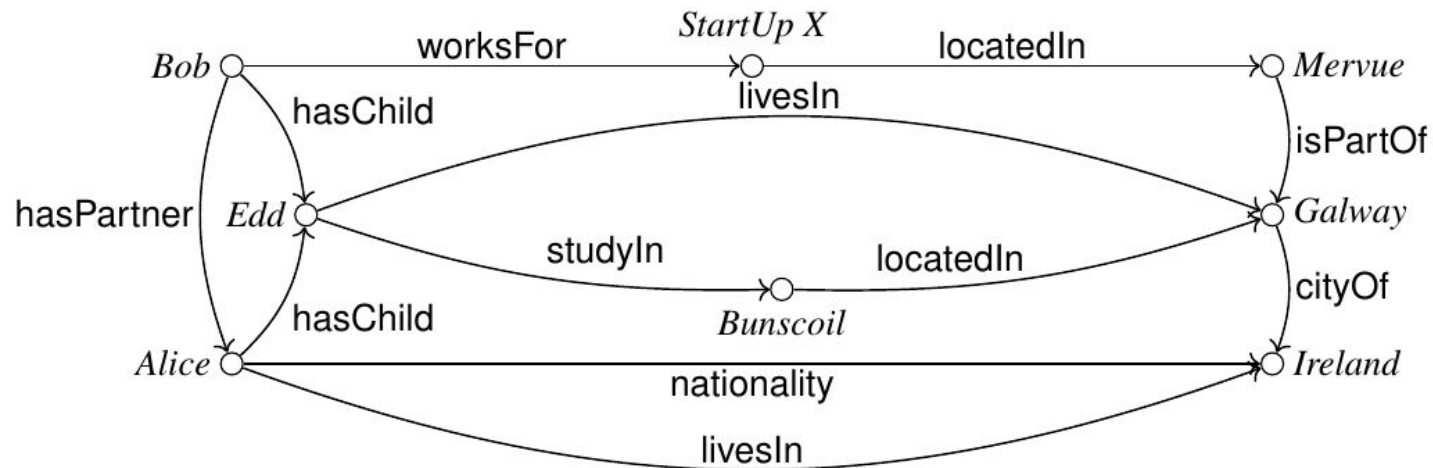


Walk path:  $\text{studyIn}^{-1}, \text{hasChild}^{-1}$

Walk path:  $\text{locatedIn}, \text{cityOf}$

## Search Approximate Equivalences of a path (3)

- (2) Extracting subgraphs of path extension instances
- We use constrained Depth-First-Search (DFS)



**Figure 1:** Example of a knowledge graph about people living in Galway, Ireland

## Search Approximate Equivalences of a path (3)

### (3) Build connecting paths between domain and range nodes

---

#### Algorithm 2 (ConnectingPaths: CONNECTING PATHS EXTRACTION)

---

**Input:**  $v_1, v_2$  nodes, depth  $d$ , knowledge graph  $\mathcal{G}$

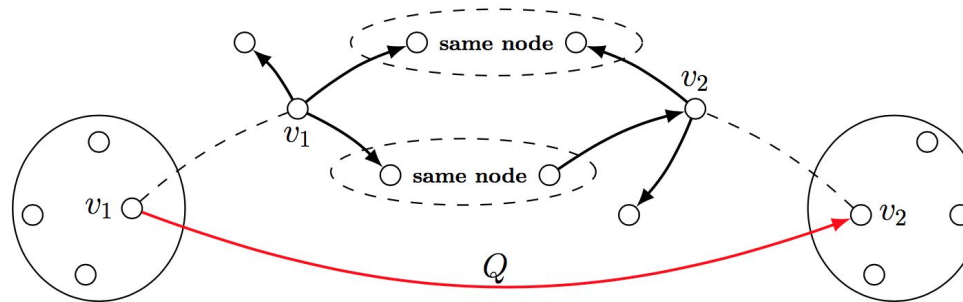
**Output:**  $\mathcal{C}$  list of connecting paths between  $v_1$  and  $v_2$

```

1:  $\mathcal{G}_1, \mathcal{G}_2 \leftarrow \text{Graph}^d(\mathcal{G}, v_1), \text{Graph}^d(\mathcal{G}, v_2)$ 
2:  $T_1, T_2 \leftarrow \{v \mid u, v \in V_{\mathcal{G}_1} \wedge \mathbf{I}_{\mathcal{G}_1}(P(u \rightsquigarrow v))\}, \{w \mid y, w \in V_{\mathcal{G}_2} \wedge \mathbf{I}_{\mathcal{G}_2}(P(y \rightsquigarrow w))\}$ 
3: for  $t \in T_1 \cap T_2$  do
4:   for  $P_1 \in \{P \mid u \in V_{\mathcal{G}_1} \wedge \mathbf{I}_{\mathcal{G}_1}(P(u \rightsquigarrow t))\}$  do
5:     for  $P_2 \in \{P \mid v \in V_{\mathcal{G}_2} \wedge \mathbf{I}_{\mathcal{G}_2}(P(v \rightsquigarrow t))\}$  do
6:        $\mathcal{C}.\text{append}(P_1 \oplus \text{Inverse}(P_2))$ 
7: return  $\mathcal{C}$ 

```

---



**Fig. 3:** Generation of connecting paths from the path extension of query  $Q$ .

## Search Approximate Equivalences of a path (4)

### (4) Ranking candidate relation paths

Rank extracted connecting paths using the following tri-criteria ranking function:

$$Rank_Q(P) = \alpha \underbrace{\frac{|PEXT_{\mathcal{G}}(P)|}{|PEXT_{\mathcal{G}}(Q)|}}_{\text{CR-1}} + \beta \underbrace{\frac{\sigma(P)}{\max\{\sigma(P_i) : P_i \in \mathcal{C}\}}}_{\text{CR-2}} + \gamma \underbrace{\frac{|Q| - |P|}{\max\{|Q|, |P|\}}}_{\text{CR3}},$$

Similarity with Query Path	Frequency Among Candidate Paths	Path Length
----------------------------	---------------------------------	-------------

where  $\alpha, \beta, \gamma$  are configurable parameters.

# Experimental setup

**GOAL:** Find ranked  $\Delta$ -equivalences for a given query path

**We experimented our approach on 4 datasets:**

- **NELL**
- **DBpedia**
- **YAGO3**
- **WordNet**

**Table 1:** Statistics of knowledge graphs used in our experiments.

	NELL	DBpedia	YAGO3	WordNet
#Entities - $ V $	1.2M	1.2M	2.6M	10K
#Relations - $\Sigma_E$	520	644	36	18
#Triples - $ E $	3.8M	4M	5.5M	141K

## Results (1)

- Queries from DBpedia dataset

Query	$Rank_Q(P)$	CR-1	CR-2	CR-3
Query B.1: $\langle \text{wasBornIn}, \text{isLocatedIn} \rangle$	Time ca. 118 min.		1000 <sup>b</sup> instances	
$\langle \text{wasBornIn}, \text{isLocatedIn}, \text{isLocatedIn}^{-1}, \text{isLocatedIn} \rangle$	1.276	0.776	1.000	-0.500
$\langle \text{isCitizenOf} \rangle$	1.025	0.016	0.001	0.500
$\langle \text{isPoliticianOf} \rangle$	0.517	0.013	0.001	0.500
$\langle \text{livesIn} \rangle$	0.514	0.007	0.000	0.500
$\langle \text{hasGender}, \text{hasGender}^{-1}, \text{isPoliticianOf} \rangle$	0.507	0.332	0.243	-0.333
Query B.2: $\langle \text{actedIn}, \text{directed}^{-1} \rangle$	Time ca. 193 min.		1000 <sup>b</sup> instances	
$\langle \text{actedIn}, \text{isLocatedIn}, \text{isLocatedIn}^{-1}, \text{directed}^{-1} \rangle$	1.360	0.860	1.000	-0.500
$\langle \text{hasGender}, \text{hasGender}^{-1} \rangle$	0.587	0.583	0.004	0.000
$\langle \text{actedIn}, \text{actedIn}^{-1}, \text{actedIn}, \text{directed}^{-1} \rangle$	0.524	0.674	0.350	-0.500
$\langle \epsilon \rangle^a$	0.518	0.018	0.000	0.500
$\langle \text{isMarriedTo} \rangle$	0.503	0.003	0.000	0.500

## Results (2)

- Queries from Yago dataset

Query	$Rank_Q(P)$	CR-1	CR-2	CR-3
Query C.1: $\langle \text{artist}, \text{bandMember} \rangle$	Time ca. 58 min.		1000 <sup>b</sup> instances	
$\langle \text{artist}, \text{associatedMusicalArtist}^{-1}, \text{associatedBand}, \text{bandMember} \rangle$	1.435	0.935	1.000	-0.500
$\langle \text{artist}, \text{associatedBand}^{-1}, \text{associatedMusicalArtist}, \text{bandMember} \rangle$	1.429	0.935	0.994	-0.500
$\langle \text{artist}, \text{associatedBand}^{-1}, \text{associatedMusicalArtist}, \text{associatedBand}^{-1} \rangle$	0.953	0.524	0.929	-0.500
$\langle \text{artist}, \text{associatedMusicalArtist}^{-1}, \text{associatedBand}, \text{associatedMusicalArtist}^{-1} \rangle$	0.952	0.524	0.928	-0.500
$\langle \text{genre}, \text{instrument}, \text{instrument}^{-1}, \text{genre}^{-1} \rangle$	0.736	0.432	0.804	-0.500
Query C.2: $\langle \text{academicAdvisor}, \text{almaMater} \rangle$	Time ca. 80 min.		335 instances	
$\langle \text{academicAdvisor}, \text{birthPlace}, \text{birthPlace}^{-1}, \text{almaMater} \rangle$	1.080	0.580	1.000	-0.500
$\langle \text{academicAdvisor}, \text{deathPlace}, \text{birthPlace}^{-1}, \text{almaMater} \rangle$	0.846	0.575	0.771	-0.500
$\langle \text{almaMater} \rangle$	0.641	0.121	0.020	0.500
$\langle \text{academicAdvisor}, \text{deathPlace}, \text{deathPlace}^{-1}, \text{almaMater} \rangle$	0.587	0.620	0.467	-0.500
$\langle \text{notableStudent}^{-1}, \text{almaMater} \rangle$	0.540	0.459	0.081	0.000



## Results (3)

- Queries from NELL dataset

Query	$Rank_Q(P)$	CR-1	CR-2	CR-3
Query A.1: $\langle \text{riverEmptiesIntoRiver}, \text{riverFlowsThroughCity} \rangle$	Time ca. 53.2 min.	519 instances		
$\langle \text{cityLiesOnRiver}^{-1}, \text{generalizations}, \text{generalizations}^{-1} \rangle$	1.342	0.688	0.987	-0.333
$\langle \text{riverFlowsThroughCity}, \text{generalizations}, \text{generalizations}^{-1} \rangle$	1.337	0.688	0.982	-0.333
$\langle \text{cityLiesOnRiver}^{-1}, \text{generalizations}, \text{generalizations}^{-1}, \text{generalizations}^{-1} \rangle$	1.192	0.692	1.000	-0.500
$\langle \text{riverFlowsThroughCity}, \text{generalizations}, \text{generalizations}^{-1}, \text{generalizations}^{-1} \rangle$	1.189	0.692	0.997	-0.500
$\langle \text{riverEmptiesIntoRiver}, \text{cityLiesOnRiver}^{-1} \rangle$	1.015	0.996	0.019	0.000
Query A.2: $\langle \text{athletePlaysForTeam}, \text{teamHomeStadium}, \text{stadiumLocatedInCity} \rangle$	Time ca. 35 min.	326 instances		
$\langle \text{athletePlaysForTeam}, \text{generalizations}, \text{generalizations}^{-1}, \text{citySportsTeams}^{-1} \rangle$	1.404	0.770	0.884	-0.250
$\langle \text{athletePlaysForTeam}, \text{generalizations}, \text{generalizations}^{-1}, \text{teamPlaysInCity} \rangle$	1.353	0.764	0.839	-0.250
$\langle \text{teamMember}^{-1}, \text{generalizations}, \text{generalizations}^{-1}, \text{citySportsTeams}^{-1} \rangle$	1.264	0.739	0.775	-0.250
$\langle \text{athletePlaysForTeam}, \text{teamPlaysAgainstTeam}^{-1}, \text{teamPlaysAgainstTeam}^{-1}, \text{citySportsTeams}^{-1} \rangle$	1.249	0.531	0.969	-0.250
$\langle \text{teamMember}^{-1}, \text{generalizations}, \text{generalizations}^{-1}, \text{teamPlaysInCity} \rangle$	1.244	0.733	0.761	-0.250



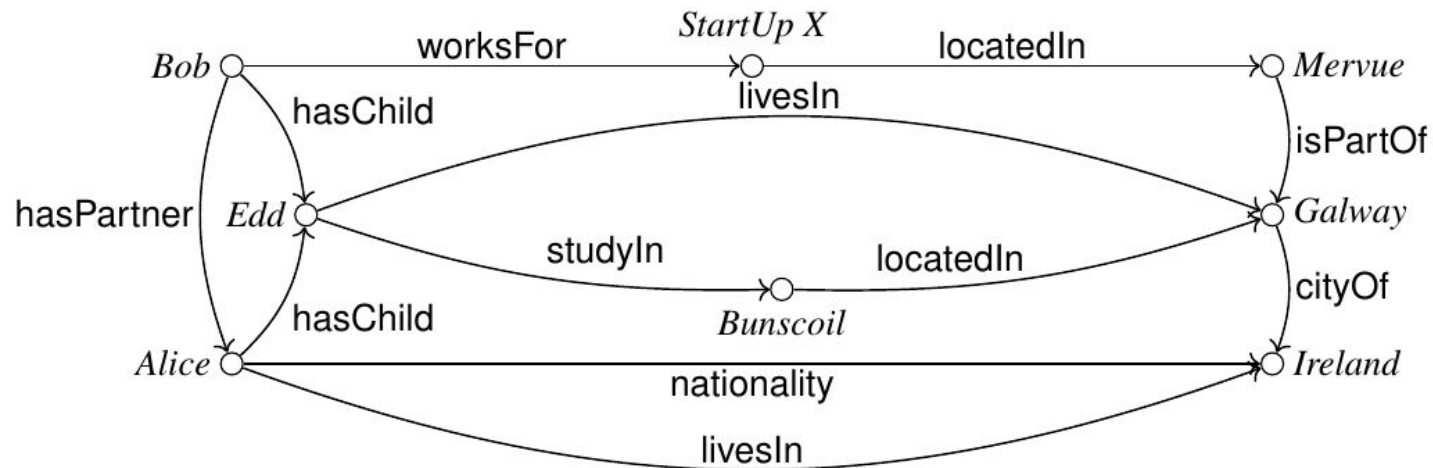
# Conclusions

Our proposed technique for identifying equivalences of a relation path that achieves the following:

- Address the complexity of finding strict equivalences by using samples of path extensions
- Query and rank approximately equivalent paths depending on multiple ranking criteria.

## Future work

1. Association rule mining by querying singular relation paths
2. Knowledge embedding to latent feature models using path similarities



**Figure 1:** Example of a knowledge graph about people living in Galway, Ireland



Insight



# Questions ?

Thank you