# Orion: Enabling Suggestions in a Visual Query Builder for Ultra-Heterogeneous Graphs

Nandish Jayaram

Member of Technical Staff 3

Pivotal

(work done as a student at the University of Texas at Arlington)

**Rohit Bhoopalam**

**Chengkai Li**
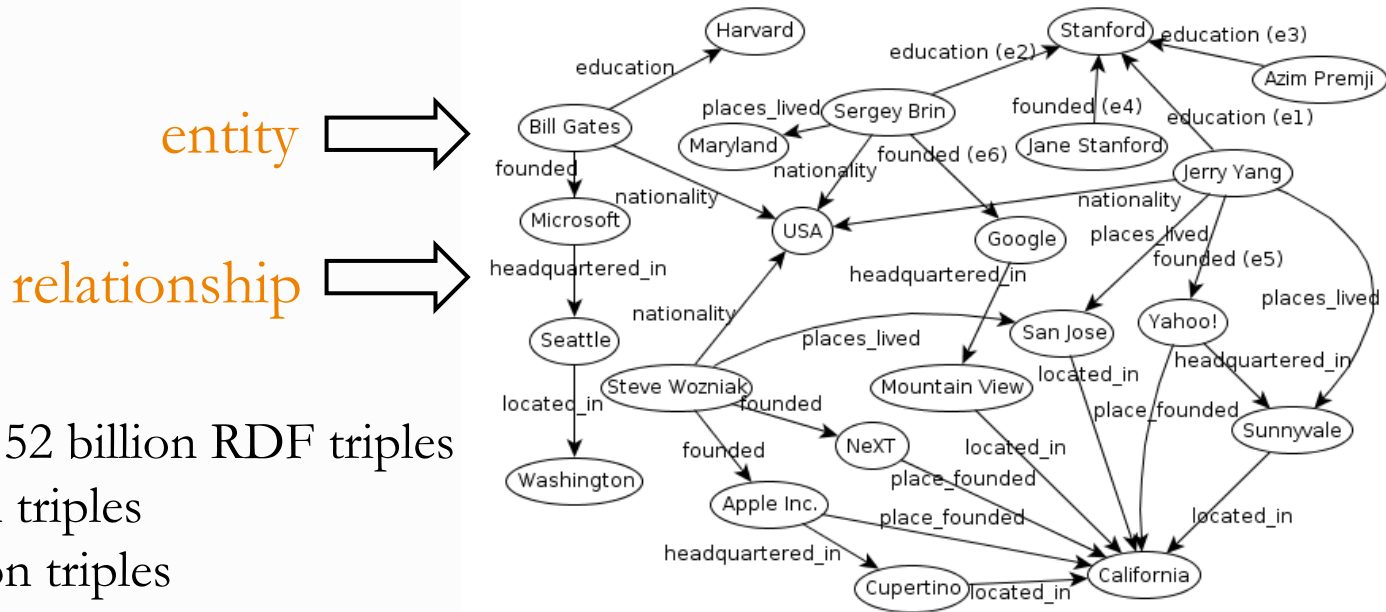
**Vassilis Athitsos**

Eighth TUC Meeting , June 22-23, 2016

# Ultra-heterogeneous Entity Graphs

Large, complex and schema-less graphs capturing millions of entities and billions of relationships between entities.

entity ⟹

relationship ⟹



Linked Open Data : 52 billion RDF triples
Freebase : 1.8 billion triples
DBpedia : 470 million triples
Yago : 120 million triples

# Structured Queries are Difficult to Write
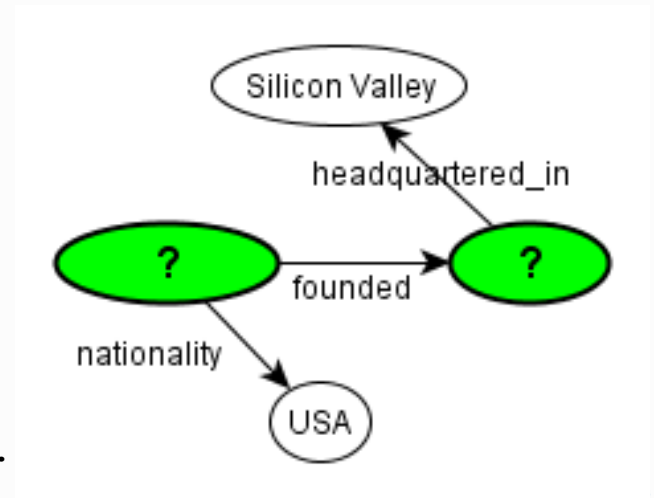
**SQL QUERY:**

```
SELECT  Founder.subj, Founder.obj
FROM Founder,
    Nationality,
    HeadquarteredIn
WHERE
    Founder.subj = Nationality.subj AND
    Founder.obj = HeadquarteredIn.subj
```

**SPARQL QUERY:**

```
SELECT  ?company  ?founder  WHERE {
    ?founder  dbo:founded  ?company .
    ?founder  dbo:nationality  db:United_States .
    ?company  dbprop:headquartered_in  db:Silicon_Valley .
}
```

- Require knowledge on data model, query language, and schema.

- Well-known usability challenges [Jagadish+07]



**3**

# Simpler Query Paradigms

## Keyword Search

o [Kargar+11], BLINKS [He+07]

  o Challenging to articulate exact query intent by keywords

## Approximate Query Answering

o NESS [Khan+11]: uses neighborhood-based indexes to quickly find approximate matches to a query graph;

o TALE [Tian+08]: approximate large graph matching

  o Users still have to formulate the initial query graph

# Visual Query Builders

**Relational Databases:** CLIDE [Petropoulos+06]

**Graph Databases:** VOGUE [Bhowmick+13], PRAGUE [Jin+12], Gblender [Jin+10], GRAPHITE [Chau+08]

**Single Large Graphs:** QUBLE [Hung+13]

o   Require a good knowledge of the underlying schema

o   No automatic suggestions regarding what to include in the query graph

# Orion

o Interactive GUI for building query components

o Iteratively suggests edges based on their relevance to the user's query intent, according to the partial query graph so far

# Orion GUI

Dynamic list of all possible user actions at any given moment

Control panel for various settings and tips

**Possible Actions**

**Click** on other grey nodes to be included in the query graph.

**Click** on the grey edge to select it, or click on a grey edge to display the other occurrences of the grey edgei, if any.

**Click** on the empty canvas to add the selected nodes and edges to the query graph while ignoring the unselected grey nodes, and display new suggestions.

**Click** on selected nodes (in blue) to unselect them.

Submit

Useful Tips    +

Edge Types    +

Settings    +

Clear Canvas

Refresh Suggestions

film_location/featured_in_films

FILM

film/produced_by

FILM PRODUCER

film/directed_by

FILM DIRECTOR

film-starring

United States of America

person/place_of_birth

FILM ACTOR

nominated_for-award_nominations

places_lived-location

LOCATION

AWARD-NOMINATED WORK

# Active Mode

Grey edges and nodes automatically suggested in **active mode:**

o   Accepted by user (blue): positive edges
o   Ignored by user: negative edges

# Passive Mode



A new node added in **passive mode**

A new edge added in **passive mode**

# Concepts

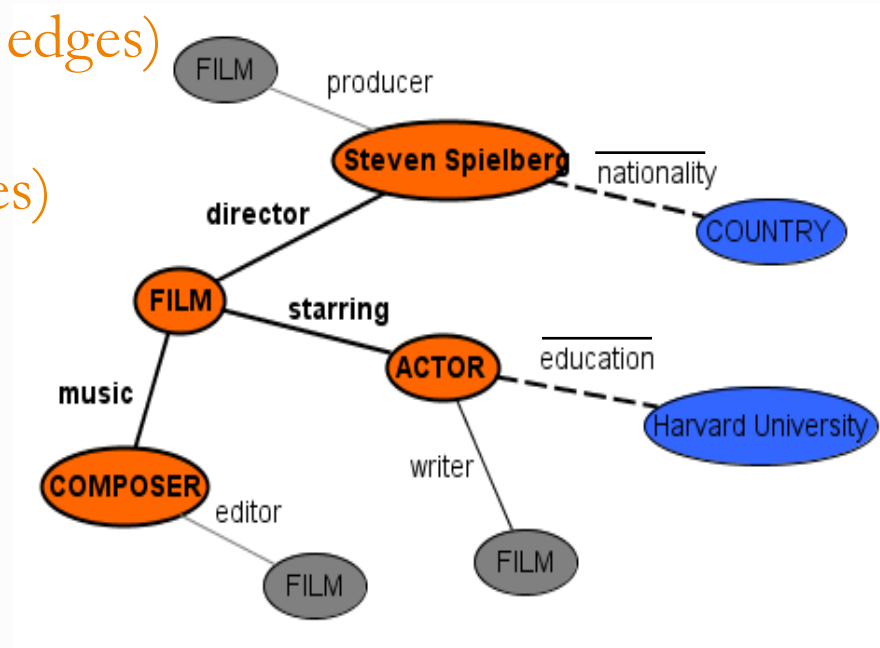Edges in partial query graph (positive edges)
    `starring, director, music`

Edges rejected by users (negative edges)
    `education, nationality`

Candidate edges
    `producer, writer, editor`



**Query Session:**
`<starring, director, music, education, nationality>`

# Concepts

## Query Log (W)

| Id | Query Session |
|---|---|
| $w_1$ | *education, founder, $\overline{nationality}$* |
| $w_2$ | *starring, $\overline{music}$, director* |
| $w_3$ | *$\underline{nationality}$, $\overline{education}$, music, $\overline{starring}$* |
| $w_4$ | *$\underline{artist}$, $\overline{title}$, writer, director* |
| $w_5$ | *$\overline{director}$, founder, producer* |
| $w_6$ | *$\underline{writer}$, $\overline{editor}$, genre* |
| $w_7$ | *$\underline{award}$, movie, director, $\overline{genre}$* |
| $w_8$ | *education, founder, $\overline{nationality}$* |

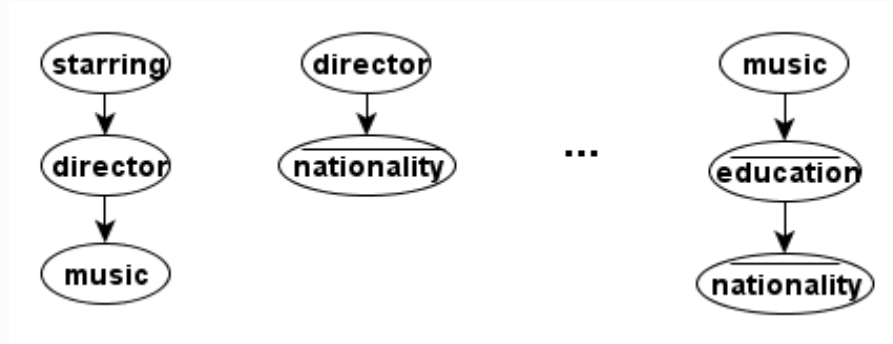Positive Edge

Negative Edge

## Problem

Given a query log, a query session, and a set of candidate edges, rank the candidate edges by their relevance to the user's query intent

# Random Decision Path (RDP)

`<starring, director, music, education, nationality>`

o Choose edges from the query session randomly, to form RDPs



o Each decision path selects a subset of the query log, with no more than 'τ' rows

o Grow a path incrementally until its support in the query log drops below 'τ'

# Random Decision Path: Scoring

o For each RDP, use its corresponding query log subset to compute the support of each candidate edge.

o Final score of each candidate is its average score across all RDPs.

o If R is the set of all RDPs:

$$score(e) = \frac{1}{|R|} * \sum_{Q_i \in R} \sup(e, Q_i, W)$$

$$\sup(e, Q_i, W) = \frac{|\{w \mid w \in W, Q_i \cup \{e\} \subseteq w\}|}{|\{w \mid w \in W, Q_i \subseteq w\}|}$$
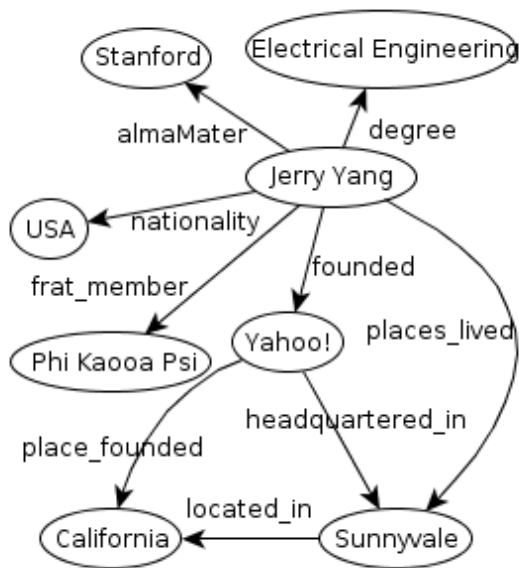
# Query Log

Nonexistent  (almost)

Simulate and bootstrap

o  Find positive edges

 o  Wikipedia and data graph

 o  Data graph only

 o  SPARQL query log [Morsey+11]

o  Inject negative edges

| Id | Query Session |
|----|---------------|
| $w_1$ | education, founder, nationancitȳ |
| $w_2$ | starring, m̄ūsīc, director |
| $w_3$ | nationality, ēducātion, music, stārrīng |
| $w_4$ | ārtīst, tītle, writer, director |
| $w_5$ | dīrēctor, founder, producer |
| $w_6$ | writer, ēditor, genre |
| $w_7$ | āward, movie, director, gēnrē |
| $w_8$ | education, founder, nationancitȳ |

# Query Log Simulation: Wikipedia + Data Graph

## Use Sentences in Wikipedia Articles to Identify Positive Edges



### Early life [edit]

Yang was born in Taipei, Taiwan on November 6, 1968, and moved to San Jose, California at the age of ten with his mother and younger brother.[4] He claimed that despite his mother being an English teacher, he only knew one English word (shoe) on his arrival. Becoming fluent in the language in three years, he was then placed into an Advanced Placement English class.[5]

Yang graduated from Sierramont Middle School and Piedmont Hills High School in San Jose and went on to earn a Bachelor of Science and a Master of Science in electrical engineering from Stanford University, where he was a member of Phi Kappa Psi fraternity.[6][7]
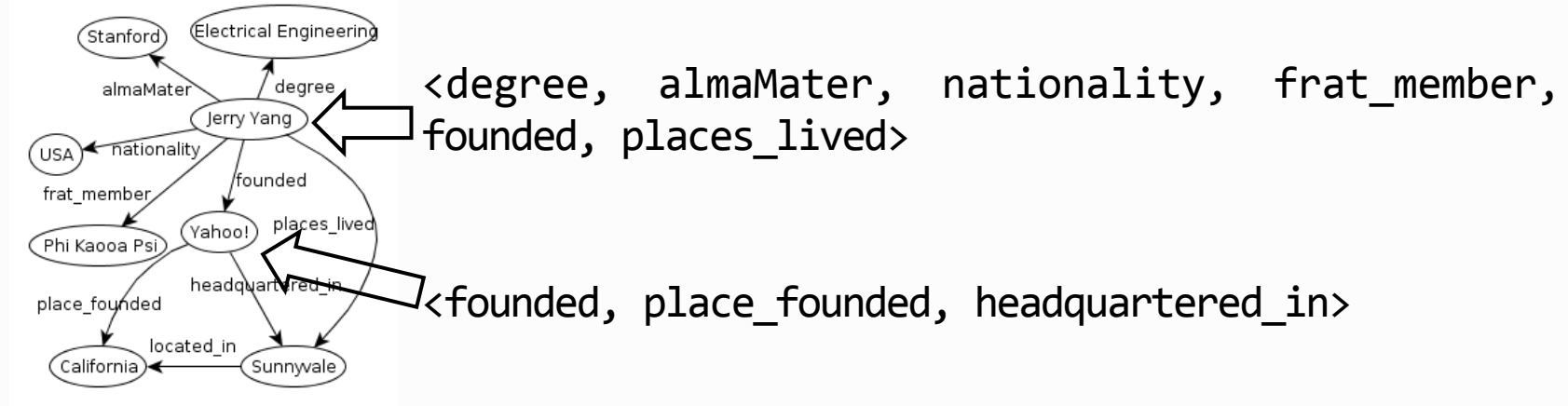
**Nodes Mapped:** **Jerry Yang,** Electrical Engineering, Stanford University, Phi Kappa Psi

`<degree, almaMater, frat_member>`

# Query Log Simulation: Data Graph Only

Represent Each Node as an Itemset of Positive Edges



`<degree, almaMater, nationality, frat_member, founded, places_lived>`

`<founded, place_founded, headquartered_in>`

Generate Frequent Itemsets of Varying Sizes

o   Each frequent itemset of edges forms positive edges

# Query Log Simulation: Injecting Negative Edges

## Positive Edges List

(1) writer, starring, producer

(2) starring, editor, education

(3) editor, nationality, music

## Inject Negative Edges

writer, starring, producer, editor, education (starring appears in 2)

starring, editor, education, writer, producer, nationality, music (starring appears in 1, and editor appears in 3)

editor, nationality, music, starring, education (editor appears in 2)

# Experiments

## System Configurations

- o Double quad-core 2.0 GHz Xeon server, 24 GB RAM
- o TACC: 5 Dell PowerEdge R910 server nodes, with 4 Intel Xeon E7540 2.0 GHz 6-core processors, 1 TB RAM

## Datasets

- o Freebase (33 M edges, 30 M nodes, 5253 edge types)
- o DBpedia (12 M edges, 4 M nodes, 647 edge types)

## User Studies with Freebase

# Query Logs Compared

o Freebase: Wiki, Data

o DBpedia: Wiki, Data, QLog

| Query Log | Components Used in Query Log Simulation | | | |
|-----------|----------|---------|-----------|-------------|
|           | Freebase | DBpedia | Wikipedia | SPARQL [26] |
| Wiki-FB   | Yes      | -       | Yes       | -           |
| Data-FB   | Yes      | -       | -         | -           |
| Wiki-DB   | -        | Yes     | Yes       | -           |
| Data-DB   | -        | Yes     | -         | -           |
| QLog-DB   | -        | -       | -         | Yes         |

# User Studies: Setup

15 Users for Orion, 15 Users for Naïve  (A/B testing)

45 Easy, 30 Medium, and 30 Hard Query Tasks Designed

3 Easy, 2 Medium, 2 Hard Queries Assigned per Query Task

105 Query Tasks per System in Total

4 Survey Questions per Query Task

| Likert Scale Score | Q1: How well do you think the query graph formulated by you captures the required query intent? | Q2: How easy was it to use the interface for formulating this query? | Q3: How satisfactory was the overall experience? | Q4: The interface provide features necessary for easily formulating query graphs. |
|---|---|---|---|---|
| 1 | Very Poorly | Very Hard | Unacceptable | Strongly Disagree |
| 2 | Poorly | Hard | Poor | Disagree |
| 3 | Adequately | Neither Easy Nor Hard | Satisfactory | Uncertain |
| 4 | Well | Easy | Good | Agree |
| 5 | Very Well | Very Easy | Excellent | Strongly Agree |

# User Studies: Conversion Rate

Conversion Rate:

- o Percentage of query tasks completed successfully
- o Successful completion measured using edge isomorphism, and not a binary notion of matching

| System | Queries | Sample Size | Conversion Rate ($c$) | z-value | p-value |
|--------|---------|-------------|----------------------|---------|---------|
| Orion | All | 105 | $c_O$=0.74 | 0.92 | 0.1788 |
| Naive | | | $c_N$=0.68 | | |
| Orion | Medium + Hard | 60 | $c_O$=0.70 | **1.36** | **0.0869** |
| Naive | | | $c_N$=0.58 | | |

Orion has a higher conversion rate than Naïve for complex queries!

# User Studies: User Experience Results



As the difficulty level of the query graph being constructed increases, the usability of Orion seems significantly better than Naïve's

# Query Logs Comparison



**Positive edges better captured based on the context of human usage of relationships in Wikipedia**

**DBpedia is created using info-boxes in Wikipedia, and is thus very clean. Wiki-DB is highly similar to Data-DB for DBpedia**

# Challenges for the LDBC Community

A benchmark query log to help improve the performance of systems such as Orion

A benchmark query set for visual query formulation, for better evaluation of systems

# Orion

## Prototype
http://idir.uta.edu/orion

## Introduction Video
http://bit.ly/1O0GnNo

# Thank You!  Questions?

# User Studies: Efficiency by Number of Iterations



Time required to construct query graphs in Orion is comparable to Naïve in most cases, despite the steeper learning curve of Orion due to more features

# User Studies: Efficiency by Time



Time required to construct query graphs in Orion is comparable to Naïve in most cases, despite the steeper learning curve of Orion due to more features

# Edge Ranking Algorithms: Efficiency by Time



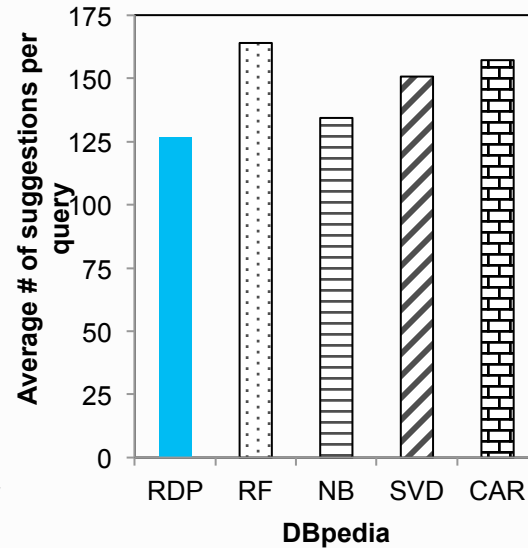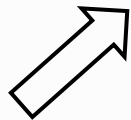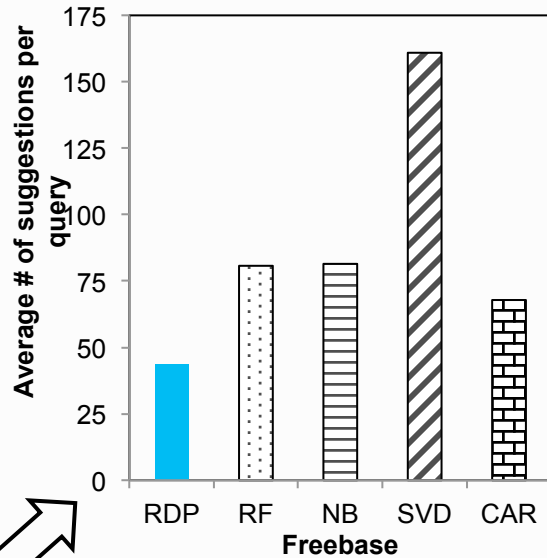RDP better than RF and comparable to NB, despite RF and NB being light models

RDP significantly better than SVD and CAR, but worse than RF and NB

# Edge Ranking Algorithms

o    Simulates only Active Mode

o    43 target query graphs for Freebase

   o 6 two-edged, 10 three-edged,  9 four-edged, 17 five-edged, 1 six-edged (includes medium and hard queries from the user study)

   o 167 input instances

o    33 target query graphs for DBpedia

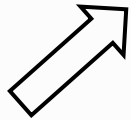   o 2 three-edged,  29 four-edged, 2 five-edged

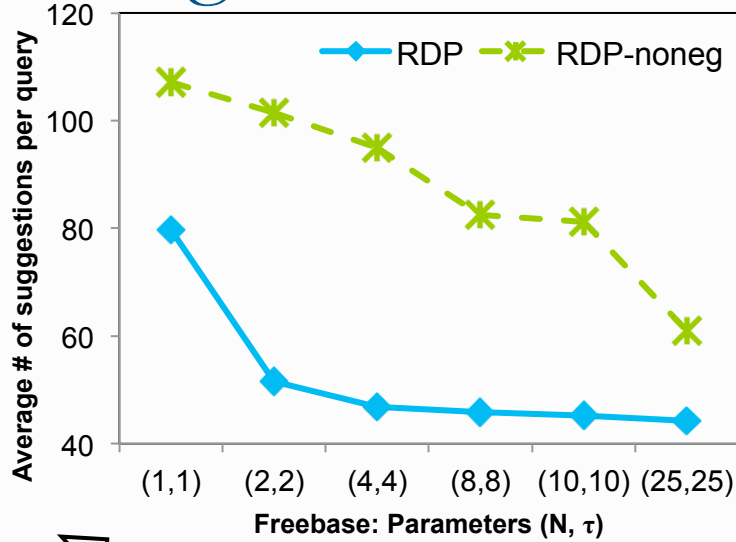   o 130 input instances

# Edge Ranking Algorithms: Efficiency by Number of Suggestions
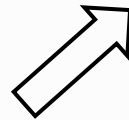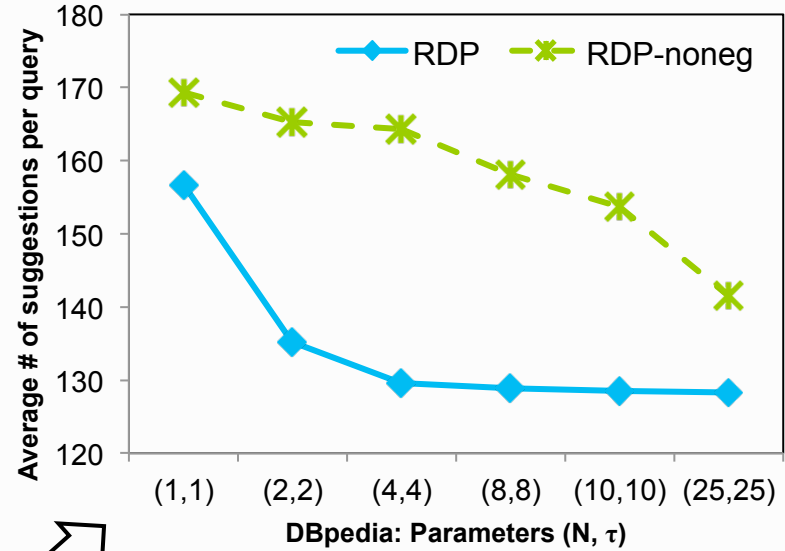


RDP requires only 40 suggestions, 1.5-4 times fewer than other methods

RDP requires fewer suggestions compared to all other methods

# Tuning RDP Parameters



RDP performs better with more random decision paths and higher query log threshold

Considering negative edges in query session is important, as it results in better performance of RDP