

# RDivF: Diversifying Keyword Search on RDF Graphs

Nikos Bikakis<sup>1,2</sup> • Giorgos Giannopoulos<sup>1,2</sup> • John Liagouris<sup>1,2</sup>  
Dimitrios Skoutas<sup>2</sup> • Theodore Dalamagas<sup>2</sup> • Timos Sellis<sup>3</sup>



<sup>1</sup> National Technical University of Athens, Greece

<sup>2</sup> ATHENA Research Center, Greece

<sup>3</sup> RMIT University, Australia



## Intro

### RDivF: RDF + Diversity

RDivF is a *diversification framework* for keyword search on RDF data. RDivF aims at exploiting several aspects of the *RDF data model* (e.g., *resource content*, *RDF graph structure*, *schema semantics*) to answer RDF keyword queries with a set of diverse results.

### Overview

#### From Keywords to Diverse RDF Subgraphs

In: Keyword query & RDF Data

Out: Ranked list of diverse RDF subgraphs

#### Retrieval Model

- Results are defined as RDF subgraphs
- RDF properties (relations) are treated as first-class citizens
- Structural and semantic homogeneity among nodes, edges and paths of the same graph result and heterogeneity between different graph results

#### Ranking Model

- Textual similarity
- Classes and properties hierarchies
- RDF/S – OWL schema semantics

## Retrieval Model

### Query Result

Assume an *RDF graph*  $G(V,E)$ , where  $V$  is the set of *vertices* and  $E$  the set of *edges*.

Let  $q = \{\{t_1, t_2, \dots, t_m\}, k, \rho\}$  be a *keyword query* comprising a set of  $m$  *terms*, a parameter  $k$  specifying the *maximum number of results* to be returned, and a parameter  $\rho$  that is used to restrict the *maximum path length* between keyword nodes.

A subgraph  $G_q$  of  $G$  is a *query result* of  $q$  over  $G$ , iff:

- for each keyword  $t$  in  $q$ , there exists exactly one node  $v$  in  $G_q$  such that  $v \in V_t$  (these are called *keyword nodes*)
- for each pair of keyword nodes  $u, v$  in  $G_q$ , there exists a path between them with length at most  $\rho$
- for each pair of keyword nodes  $u, v$  in  $G_q$ , there exists at most one direct keyword path between them
- each non-keyword node lies on a path connecting keyword nodes

### Diversified Result Set

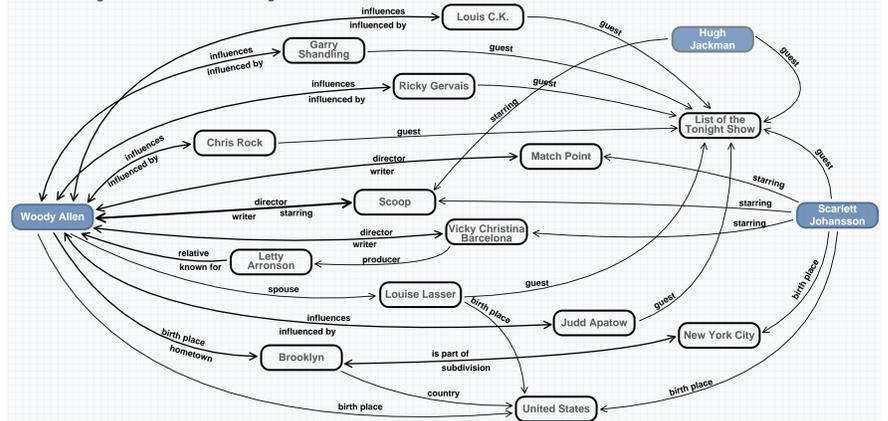
Let a function  $r: (G_q, q) \rightarrow [0, 1]$  that measures the *relevance* between the query  $q$  and a result  $G_q$ .

Let function  $d: (G_q, G'_q) \rightarrow [0, 1]$  that measures the *dissimilarity* between two query results  $G_q$  and  $G'_q$ .

The *diversified result set*  $R_k$  is a subset of the results  $R$  with size  $k$  that maximizes a combined measure of  $r$  and  $d$ .

## Motivating Example

### DBpedia Snapshot



### Example 1.

Keywords: *Woody Allen, Scarlett Johansson*

### Intermediate Results

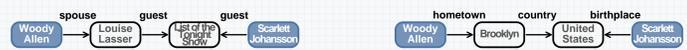
(a) Diverse paths w.r.t. properties of the searched entities



(b) Diverse paths w.r.t. the intermediate entities connecting searched entities



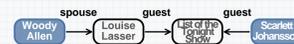
(c) Diverse paths w.r.t. the type (classes) of intermediate entities and the properties connecting searched entities



(d) Diverse paths w.r.t. the number of intermediate entities connecting searched entities



(e) Combination of c) and d)



(f) Non-diverse paths



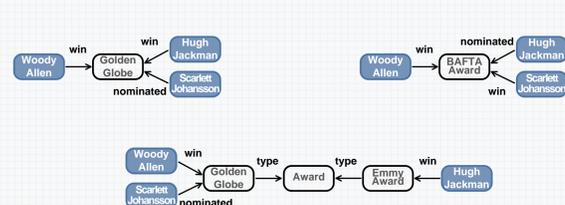
### Example 2.

Keywords: *Woody Allen, Scarlett Johansson, Hugh Jackman*

### Diverse Graph Results

The specific diversification needs should determine the ranking model and the composition of the final results

(a) Return diverse results for the searched entities, w.r.t. awards they have won



(b) Return diverse results for the searched entities, w.r.t. the people they are related to



(c) Return diverse results for the searched entities

