

A Hybrid Approach to Linked Data Query Processing with Time Constraints

Steven Lynden, Isao Kojima, Akiyoshi Matono, Akihito
Nakamura, Makoto Yui

National Institute of Advanced Industrial Science and
Technology, Japan



Motivation

- Indexing systems, e.g. Sindice, can be used to query the Semantic Web, however:
 - *Hybrid SPARQL queries: fresh vs. fast results - Umbrich et al.*
 - Coherence
 - A significant proportion of data from Sindice etc. may not be up-to-date with sources.
- Existing distributed SPARQL query processing systems are often very unpredictable in terms of response time.
- Some applications may require a best effort in a fixed amount of time
 - e.g. a portal for browsing a Linked Data repository attempting to suggest related RDF data from other sources requiring answers from a query processing back-end within the average time a user stays on a page.

Proposed approach

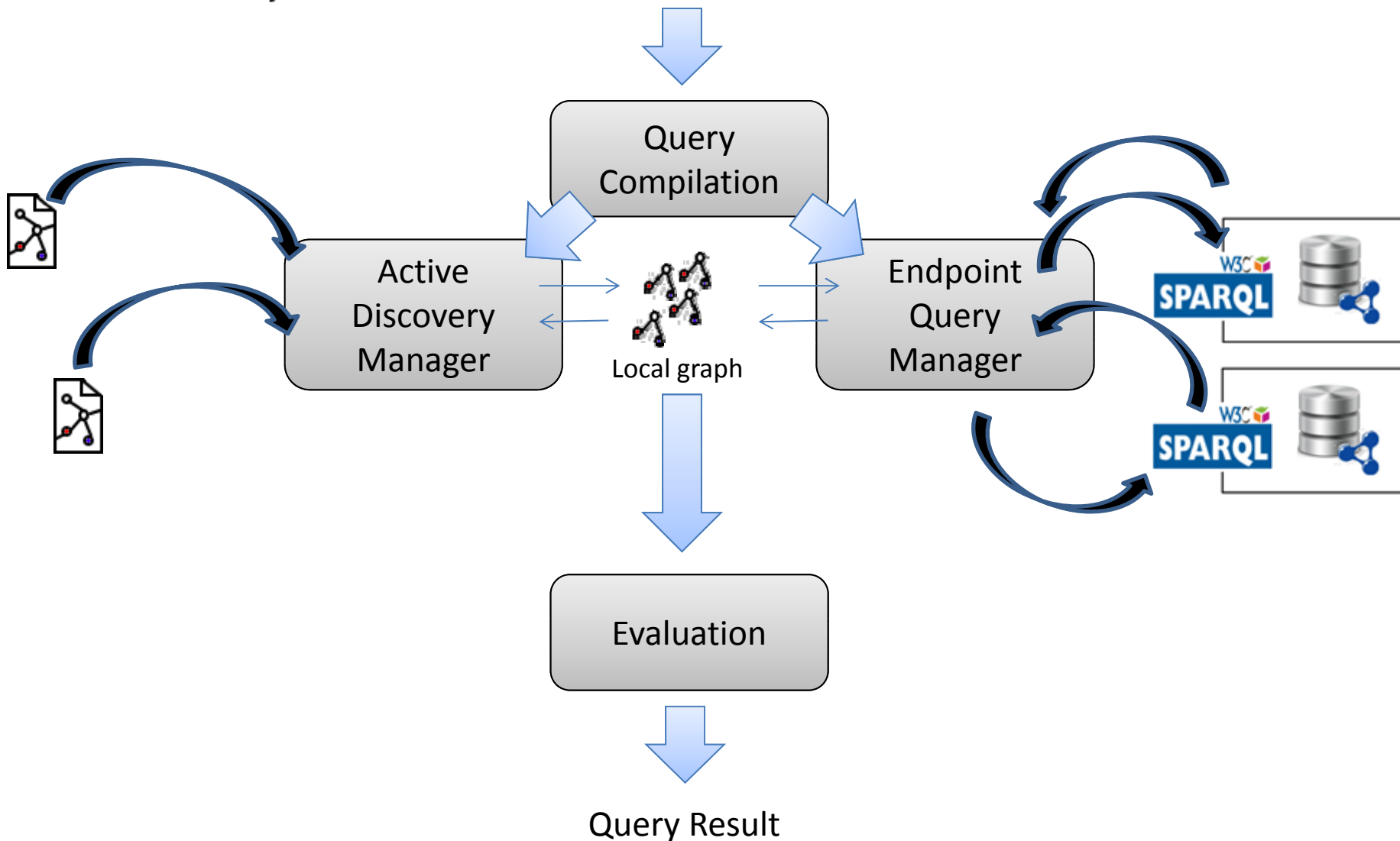
- Execute two components in parallel
 - Active discovery
 - Investigate URIs, retrieve RDF data, match against triple patterns in the query applying FILTER predicates
 - Query SPARQL endpoints
 - Construct sub-queries from the federated query, execute them using available SPARQL endpoints
- Both components share a local graph data structure in which a temporary result is constructed
- After a set time period, both components terminated and the local graph transformed into a query result

Hybrid Query Processing with Time Constraints

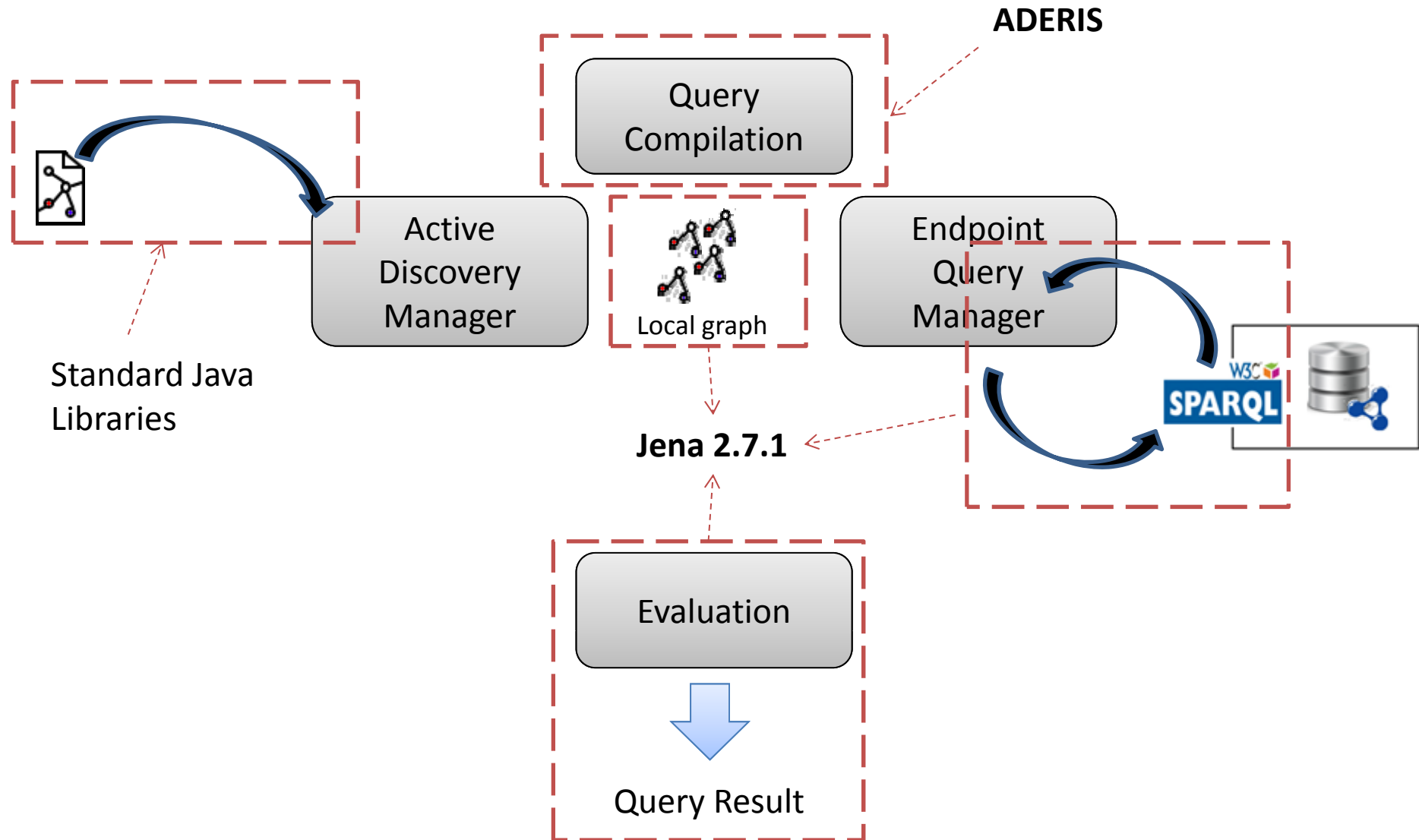
- Compile Query
- Access SPARQL endpoints and documents containing RDF data
- Stop and evaluate

User's SPARQL Query

```
PREFIX dc: <http://purl.org/dc/elements/1.1/>
PREFIX dbp: <http://dbpedia.org/resource/property>
SELECT * WHERE {
  ?x dc:subject dbp:FIFA_World_Cup-winning_countries .
  ?p dbp:managerclubs ?x .
  ?p foaf:name "Luiz Felipe Scolari"@en .
}
```



Implementation



Endpoint Query Manager

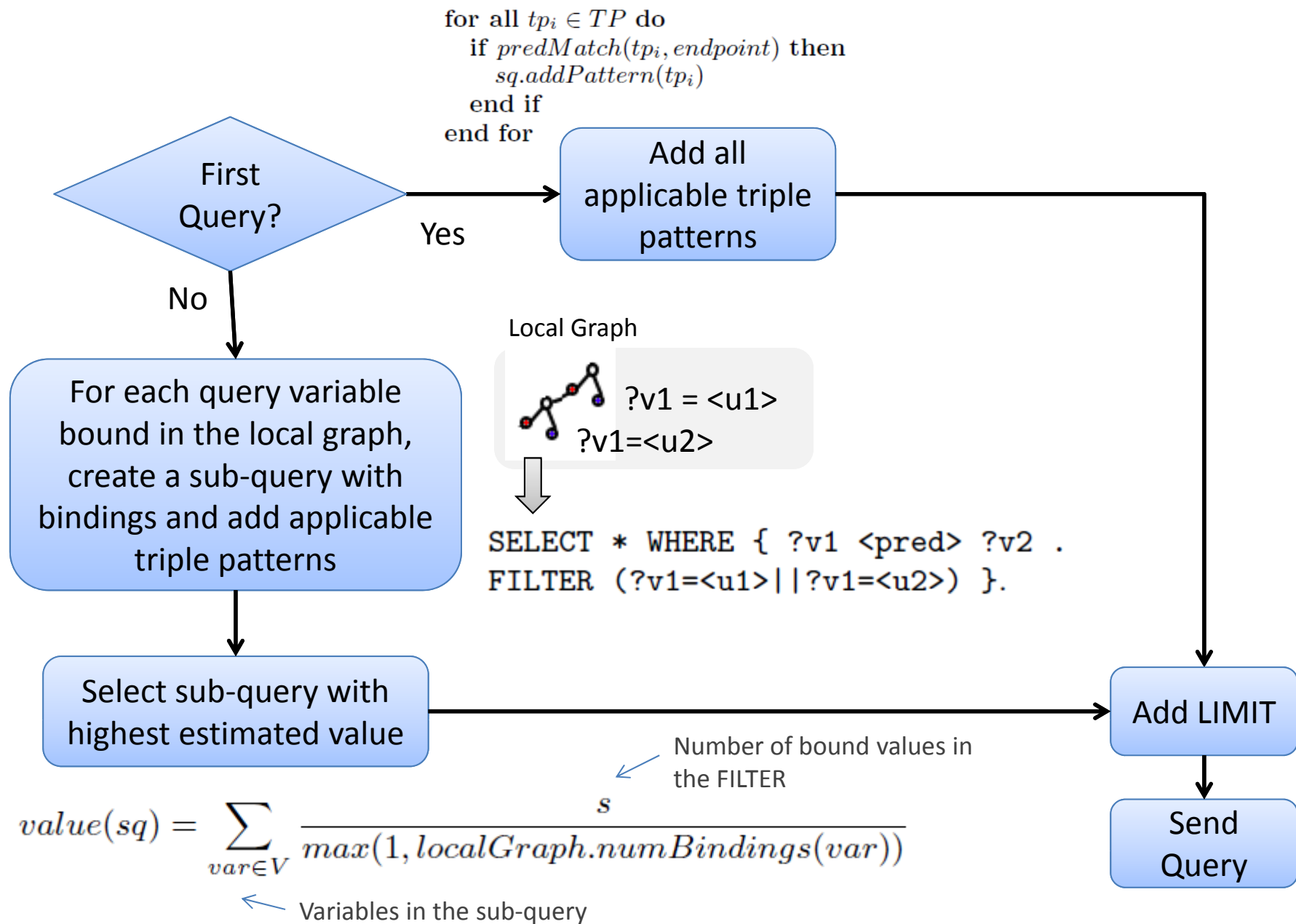
- Prior to query execution the system is configured with a set of endpoints to be used
- Existence of triples with a given predicate assumed to be known, e.g:

?paper <<http://swrc.ontoware.org/ontology#author>> ?p

triple pattern matches exist in the **data.semanticweb.org** endpoint

(Predicates in query triple patterns are usually not variables)

- Objectives
 - Execute simple queries to provide results quickly that can be explored by the active discovery manager in parallel
 - Avoid placing excessive burden on endpoints and avoid fair-use restrictions



Active Discovery Manager

- The active discover manager starts a thread for each Pay Level Domain (PLD) present in URIs in the query and as they are added to the local graph.
- Each thread is able to choose two URIs to investigate each second.
- Objective:
 - Match triple patterns in the query with RDF data retrieved via dereferencing the URIs

DBpedia URIs investigated and the number of triples matching triple patterns in the query.

```

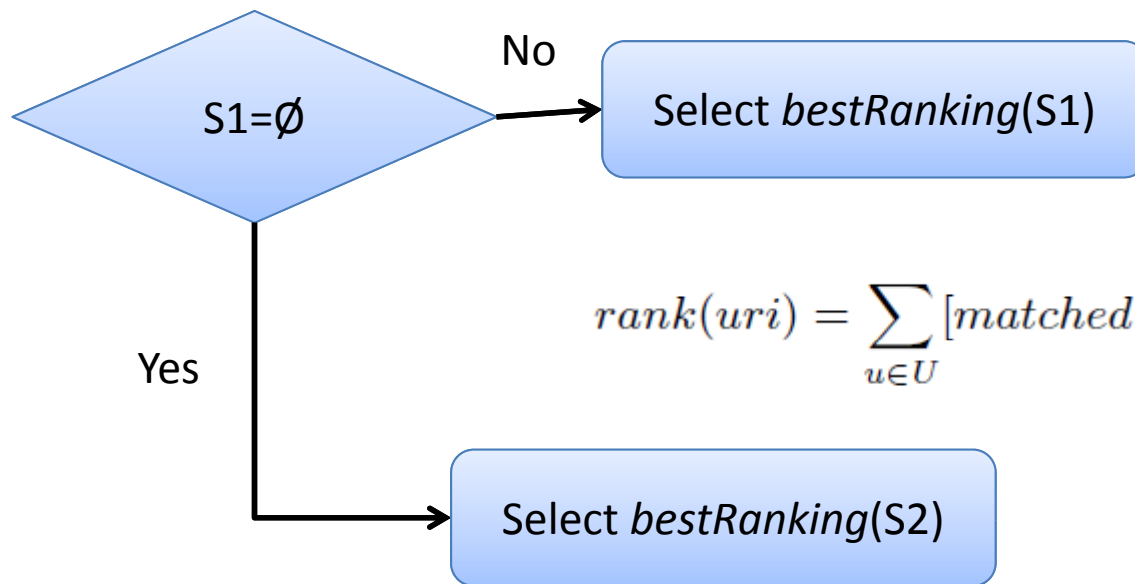
PREFIX dc: <http://purl.org/dc/elements/1.1/>
PREFIX dbp: <http://dbpedia.org/resource/property>
SELECT * WHERE {
  ?x dc:subject dbp:FIFA_World_Cup-winning_countries .
  ?p dbp:managerclubs ?x .
  ?p foaf:name "Luiz Felipe Scolari"@en .
}

```

```

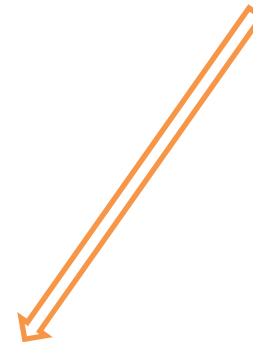
/resource/England_national_football_team (25)
/resource/Spain_national_football_team (23)
/resource/Brazil_national_team (18)
/resource/FC_Bunyodkor (3)
/resource/Vicente_Feola (2)
/resource/Luiz_Felipe_Scolari (16)

```



$$rank(uri) = \sum_{u \in U} [matched(u) * (1 - distance(u, uri))]$$

Levenshtein distance



Evaluation

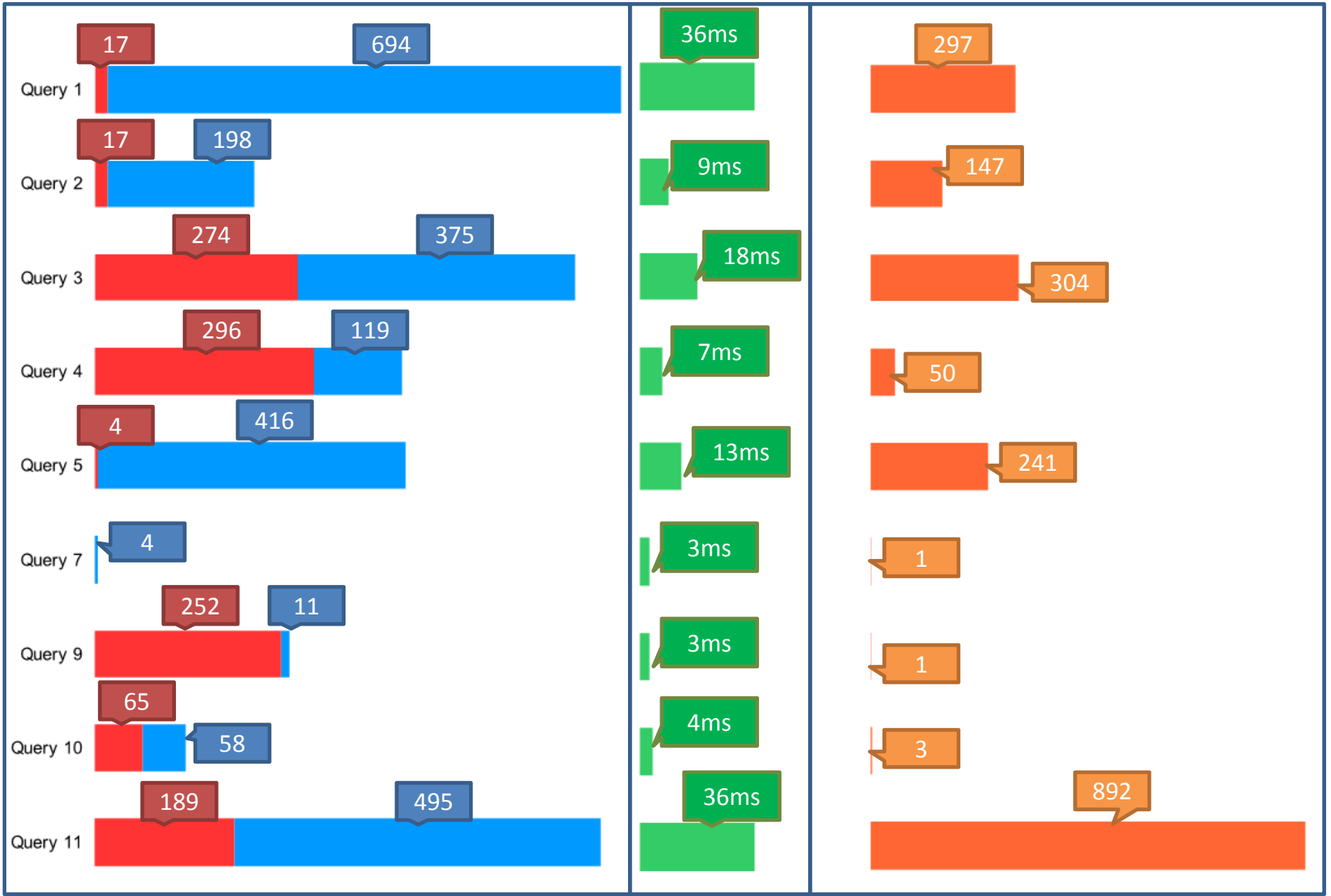
- **FedBench**
 - Benchmark for testing the efficiency and effectiveness of federated query processing on semantic data.
- Multiple query sets, we used the Linked Data (LD) query set.
- 11 Queries, however some problems encountered with 2 of the queries.
- Remaining queries executed using the proposed approach with a limit of 10 seconds.

Triples retrieved

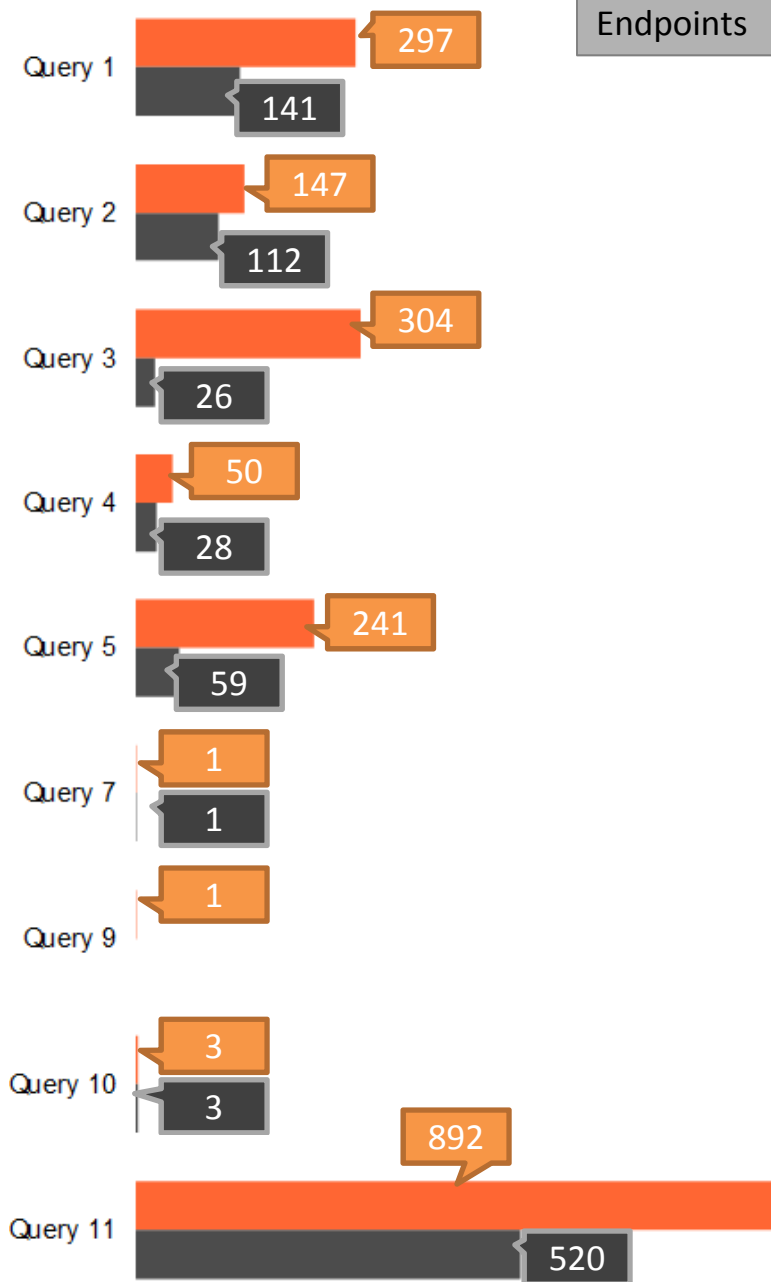
Active Discovery
SPARQL Endpoints

Eval time

results



Hybrid
ADM only (10 mins)



Endpoints

	Sindice	DBpedia	sw.org
Query 1	✓		✓
Query 2	✓		✓
Query 3	✓		✓
Query 4	✓		✓
Query 5		✓	✓
Query 7			✓
Query 9		✓	✓
Query 10			✓
Query 11			✓

PLDs

- 6
- 2
- 8
- semanticweb.org (1)
- dbpedia.org (1)
- geonames.org (1)
- dbpedia.org (1)
- 5
- dbpedia.org (1)

ADM sources with last modified < 24hrs



Conclusions

- Answering the FedBench Linked Data queries in accordance with our objective of within 10 seconds was possible using the proposed technique.
- Advantages include:
 - Fault tolerance
 - Freshness
 - Increased coverage
 - Mitigation of fair-use restrictions
- Future work will investigate benefits with more dynamic data, e.g. RDFa etc and optimisation based on relevance /quality of data sources