29 May 2011

### **RDFPath** Path Query Processing on Large RDF Graphs with MapReduce

1<sup>st</sup> Workshop on High-Performance Computing for the Semantic Web (HPCSW 2011)

Martin Przyjaciel-Zablocki Alexander Schätzle Thomas Hornung Georg Lausen

University of Freiburg Databases & Information Systems



### Overview

- 1. Motivation
- 2. MapReduce
- 3. RDFPath
- 4. Experimental Results
- 5. Summary



### Large RDF Graphs

#### Semantic Web

- Amount of Semantic Data increases steadily
- RDF is W3C standard for representing Semantic Data

#### Social Networks

- > 500 million active users in Facebook
- Interacting with >900 million objects (pages, groups, events)
- Social Graphs can also be represented in RDF
- How to explore, analyze such large RDF Graphs?
- Our Approach: Distributed analysis of large RDF Graphs by mapping RDFPath to MapReduce



### MapReduce

#### MapReduce

- Popularized by Google & widely used
- Automatic parallelization of computations
- Fixed two-stage process:  $Map \rightarrow Reduce \rightarrow Map \dots$

### Distributed File System

- Clusters of commodity hardware
   → Fault tolerance by replication
- Very large files / write-once, read-many pattern

### Apache Hadoop

- Well-known Open-Source implementation
- Used by Yahoo, Facebook, Amazon, IBM, Last.fm, ...



### MapReduce (2)

#### Steps of a MapReduce execution



#### Signatures

Map: (in\_key, in\_value) → *list* (out\_key, intermediate\_value)
 Reduce: (out\_key, *list* (intermediate\_value)) → *list* (out\_value)

#### 2. MapReduce

BURG

# 3. RDFPath

#### >>> Path queries on RDF Graphs



RDFPath: Path Query Processing on Large RDF Graphs 6

### RDFPath

### Idea

- Navigational queries on RDF Graphs
- Declarative path specification inspired by XPath
- Particularly designed with regard to MapReduce
   > Every location step is mapped to one MapReduce job
- Extendibility

### Features

- Fixed-length paths, shortest paths
- Aggregate functions, degree of a node
- Adjacent nodes and edges





### Example (1)

> Allen :: knows [country=equals("DE")] > age

- Allen (knows) Jacob [country="DE"] (age) 42
- → Allen (knows) Sarah [country="CH"]
- Allen (knows) Chris [country="CH"]



examples features

## Example (2)

Allen :: knows(\*3)

- Allen (knows) Jacob
- Allen (knows) Jacob (knows) Emily
- Allen (knows) Chris
- Allen (knows) Sarah



### Implementation

System Architecture



#### Storing RDF Graphs

- RDF Graph is loaded in advance
- Vertical Partitioning by predicate
- Sequence Files stored in HDFS contain PathObjects



### Implementation (2)

### Query Processing



3. RDFPath Implementation

BURG



### Implementation (3)

#### Mapping to MapReduce Jobs



#### Map task:

- Tagging intermediate paths and country partition for join
- Applying filter condition
- Reduce task:
  - Perform join and store resulting paths back to HDFS



3. RDFPath Implementation

# 4. Experimental Results

### >>> Properties of RDFPath



RDFPath: Path Query Processing on Large RDF Graphs 13

### Experiments

#### Environment Setup

- Cluster of 10 machines (Dual Core CPU, 4GB RAM, 1GB HD)
- Cloudera's Distribution for Hadoop 3 Beta 1
- Default configuration with 9 reducers (one per HD)
- Real–World Last.fm dataset
  - 225 million RDF triples
- Different input-sizes instead of varying the number of machines



### Query 1

The Disco Boys - I Surrender :: trackSimilar(\*4) > trackAlbum



- Linear scaling behavior for execution times & amount of data
- Execution times mainly influenced by number of intermediate results stored locally as well as transferred data

4. Evaluation Query 1

BURG

#### <u>queries</u>

### Query 3

• Chris :: knows(\*X),  $1 \le X \le 10$ 



98% of all reachable persons can be reached by at most 7 edges
 → Corresponds to the well-known six degrees of separation paradigm

Linear scaling behavior, mainly determined by number of joins



#### <u>queries</u>

### Query 3

• Chris :: knows(\*X),  $1 \le X \le 10$ 



98% of all reachable persons can be reached by at most 7 edges
 → Corresponds to the well-known six degrees of separation paradigm

Linear scaling behavior, mainly determined by number of joins



### Summary

Amount of Semantic Web data is growing constantly!

#### RDFPath

- Intuitive syntax for path queries
- Expression of interesting graph issues such as Friend of a Friend queries
- Investigating small world properties like six degrees of separation
- Designed with MapReduce in mind

#### Execution

- Direct mapping to MapReduce
- Scaling properties of MapReduce enable analysis of large RDF Graphs

#### Evaluation

- Based on real-world data from Last.fm
- Shows linear scaling behavior in the size of the graph
- Confirms the effectiveness of our approach



5. Summary

URG



### Thanks for your attention.



RDFPath: Path Query Processing on Large RDF Graphs 19

# **Backup Slides**

**))** a) <u>Last.fm</u>

UNI FREIBURG

Thanks

- b) Examples cont.
- c) Evaluation cont.
- d) <u>RDFPath Features</u>
- e) <u>Mapping</u>
- f) <u>Comparison</u>
- g) <u>Reduce-Side-Join</u>

21

### a) Last.fm



 $\overline{\phantom{a}}$ 

### b) Example (3)

# > Allen :: knows > knows [country=prefix("C")] > age [min(18)][max(80)] .avg()

- Result
  - 34



### b) Example (4)

\* :: knows [equals("Sarah")]

- Allen (knows) Sarah
- Chris (knows) Sarah
- → Jacob (knows) Emily
- → Allen (knows) Chris



### b) Example (5)

\* :: knows > knows [equals("Sarah")]

- Allen (knows) Chris (knows) Sarah
- Allen (knows) Jacob (knows) Emily



### b) Example (6)

\* :: \* > \* [equals("Sarah")]

- Allen (knows) Chris (knows) Sarah
- Allen (knows) Jacob (knows) Emily





## b) Example (7)

- > Allen :: \* .count()
- Result

• 3



### b) Example (Last.fm)

Michael\_Jackson :: artistTracks

[trackAlbum = equals(Michael\_Jackson\_-\_Thriller)]

> trackSimilar [trackDuration = min(50000)]

> trackTopFans [userCountry = equals(DE)].

- Results
  - Michael\_Jackson (artistTracks)
     Michael\_Jackson\_-\_Beat\_It (trackSimilar)
     Michael\_Jackson\_-\_Billie\_Jean (trackTopFans) Mark
  - Michael\_Jackson (artistTracks)
     Michael\_Jackson\_-\_Someone\_in\_the\_Dark (trackSimilar)
     Rihanna\_-\_Russian\_Roulette (trackTopFans) Megan

3. RDFPath Example Last.fm

URG

## c) Query 2 (Backup)

#### Michael\_Jackson :: artistTracks

- [ trackAlbum = equals('Michael\_Jackson\_-\_Thriller') ]
- > trackSimilar [ trackDuration = min(50000) ]
- > trackTopFans [ country = equals('DE') ]





### d) RDFPath Features by Example

#### Startnode

- Chris :: knows
- \* :: knows

#### Location Step

- Chris :: knows > knows > age
- Chris :: knows (2) > age
- Chris :: \*
- Filters: equals(), prefix(), suffix(), min(), max()
  - Chris :: knows > age [min(18)] [max(67)] (6)
  - Chris :: \* > \* [equals('Peter')] (7)
  - o Chris :: knows [age = min(30)] [country = prefix('D')] > name
- Bounded Search
  - Chris :: knows [country = equals('DE')] (\*3)
- Bounded Shortest Path
  - Chris :: knows (\*3).distance('Peter')
- Aggregation Functions: count(), sum(), avg(), min() and max()
  - o Chris :: \*.count()
  - o Chris :: knows > age.avg()

SURG

# e) Mapping to MapReduce Jobs

#### Composition

- Queries in RDFPath are composed of location steps
- One location step is mapped to a join in MapReduce

### Map tasks:

- Tagging for RSJs
- Applying filter conditions

### Reduce tasks:

- Performing RSJs
- Aggregating values
- Computing Shortest paths
- Checking for Cycles



3. RDFPath Implementation

# RDF Query Languages

| Property                 | $\mathbf{RQL}$ | $egin{array}{c} {f SeRQL}, \\ {f RDQL}^3, \\ {f Triple} \end{array}$ | N3       | Versa | RxPath | $\operatorname{RPL}$ | $\frac{\rm SPARQL}{(1.0)}$ | RDFPath      |
|--------------------------|----------------|--|----------|-------|--------|----------------------|----------------------------|--------------|
| Adjacent nodes           | ±              | ±  | ±        | ±     | ×      |                      | $\checkmark$               | $\checkmark$ |
| Adjacent edges           | ±              | ±  | ×        | ×     | ×      | ×                    | $\checkmark$               | $\checkmark$ |
| Degree of a node         | ±              | ×  | $\times$ | ×     | ×      | ×                    | ×                          | $\checkmark$ |
| Path                     | ×              | ×  | $\times$ | ×     | ±      | $\pm$                | ×                          | ±            |
| Fixed-length Path        | ±              | ±  | ±        | ×     | ±      | ±                    | $\checkmark$               | $\checkmark$ |
| Distance between 2 nodes | ×              | ×  | $\times$ | ×     | ×      | ×                    | ×                          | ±            |
| Diameter                 | ×              | ×  | $\times$ | ×     | ×      | ×                    | ×                          | ×            |
| Shortest Paths           | ×              | ×  | $\times$ | ×     | ×      | ×                    | ×                          | ±            |
| Aggregate functions      | ±              | ×  | ×        | ×     | ±      | ×                    | ×                          | ±            |

(×: not supported,  $\pm$ : partially supported,  $\checkmark$ : fully supported)

**Table 1.** Comparison of RDF Query Languages (adapted from [4, 5])



f)

# f) RDF Query Languages (2)

- Adjacent nodes:
  - All nodes adjacent to a startnode

#### Adjacent edges:

- All predicates of statements involving a given startnode
- Degree of a node:
  - Number of predicates involving a startnode
- Path:
  - Find paths between a fixed startnode and an endnote of arbitrary length
- Fixed-length paths:
  - Find all paths of length 2 between startnote and endnote

#### Distance between two resources:

- (length of shortest path without a max length!)
- Diameter of a graph:
  - Diameter of a given RDF Graph.



## f) Comparison to SPARQL 1.0

### Fixed-length paths

| SELECT  | ?tmp1, | , ?tmp2, | ?tmp3, |
|---------|--------|----------|--------|
|         | ?tmp4, | ,?tmp5   |        |
| WHERE { |        |          |        |
| Allen   | knows  | ?tmp1 .  |        |
| ?tmp1   | knows  | ?tmp2 .  |        |
| ?tmp2   | knows  | ?tmp3 .  |        |
| ?tmp3   | knows  | ?tmp4 .  |        |
| ?tmp5   | knows  | ?tmp5    |        |
| }       |        |          |        |

Allen :: knows(5)

#### SPARQL

#### RDFPath



# f) Comparison to SPARQL 1.0

#### Filters

```
SELECT ?tmp1, ?tmp2, ?tmp3, ?tmp4
WHERE {
   Allen knows ?tmp1 .
   ?tmp1 age ?age1 .
   Filter (?age1 >= 20)
   ?tmp1 knows ?tmp2 .
   ?tmp2 country "DE" .
}
```

#### SPARQL

3. RDFPath

Comparison

BURG

#### RDFPath

# f) Comparison to SPARQL 1.0

### Drawbacks of SPARQL

- Limited navigational capabilities
- No Shortest Path Queries
- No Aggregate functions (count, min, max , ...)
- No Abbreviations for following same edges

### Advantages of SPARQL

- Expressive general purpose query language for RDF
- Intuitive Syntax, related to SQL
- SPARQL 1.1 adds support for some path queries



### g) Reduce–Side Join

#### • Example:





## g) Reduce–Side Join (2)

#### **Mapper Input**

#### **Mapper Output**



# g) Reduce–Side Join (3)

#### Sorting phase:

- (1) Partition according to the first keypair % #reducer
- (2) Sort within a partition according to the whole keypair

#### Consequences

- A reducer gets all values with the same first keypair
- The values within a partition contain at first all new nodes and thereafter all intermediate paths



## g) Reduce–Side Join (4)

#### **Reducer Input**

### **Reducer** Output

| Chris<br>Peter<br>Manu<br>Peter (knows) Frank (knows) Chris<br>Frank (knows) Chris<br> | Peter (knows) Frank (knows) Chris (knows) Peter<br>Peter (knows) Frank (knows) Chris (knows) Manu<br>Frank (knows) Chris (knows) Peter<br>Frank (knows) Chris (knows) Manu<br>                           |
|--|--|
| Johan  | Peter (knows) Simon (knows) Johan (knows) Frank<br>Peter (knows) Simon (knows) Johan (knows) Lukas<br>Klaus (knows) Simon (knows) Johan (knows) Frank<br>Klaus (knows) Simon (knows) Johan (knows) Lukas |
| Lukas<br>Peter (knows) Simon (knows) Johan<br>Klaus (knows) Simon (knows) Johan<br>    |  |

3. RDFPath Reduce-Side-Join

BURG

RDFPath: Path Query Processing on Large RDF Graphs 39