





# Distributed Query Processing over a Wide-Area Network

**Abdoul Macina\*, Johan Montagnat**<sup>+</sup> & Olivier Corby\* \*Université Côte d'Azur, Inria, CNRS, I3S, <sup>+</sup>Université Côte d'Azur, CNRS, I3S contact: macina@i3s.unice.fr, johan.montagnat@cnrs.fr, olivier.corby@inria.fr

## POSITIONING

### CONTEXT

• Growing volume of scientific data

 Distribution of data in many acquisition sites

Data publication for cross-analysis or secondary reuse

#### CHALLENGES

#### **Distributed Query semantics**

• To define and enforce a clear SPARQL DQP semantics

#### **Performance and reliability**

• To transform queries and generate execution plans that can be executed efficiently in parallel

### **OBJECTIVES**

- To optimize performance and reliability of Distributed querying
- To design a SPARQL-compliant parallel query engine (KGRAM-DQP)

• To demonstrate the relevance of DQP techniques in realistic use cases

### **APPROACH: Distributed Query Processing (DQP)**

### **DQP SEMANTICS OVER BOTH VERTICAL AND HORIZONTAL PARTITIONS**



#### SELECT ?team ?group ?name ?members WHERE {

?team ns:team "SPARKS".

?team ns:group ?group.

## **BGP**<sub>1</sub> from $S_1 U S_2$

?group ns:name ?name.

### ?group ns:members ?members. **BGP**<sub>2</sub> from S<sub>3</sub>

**Approaches:** 

Triple Pattern-based query processing

BGP-based query processing

Our approach: Hybrid rewriting based on predicates distribution

• To write most efficient BGPs

• BGPs for both horizontal and vertical partitions

**BGP-based for all local BGPs** 

Triple Pattern-based for distributed BGP

#### Cost Estimation

For Triple Pattern : number of instances Statistics retrieved during sources selection step : SELECT COUNT(\*) queries and others heuristics
For BGPs: *min(cost(TP<sub>i</sub>) in BGP)*Subqueries sorting
Let *Q* the initial SPARQL query
After query rewriting: *Q* is a set of *BGPs*Recursive algorithm:
Get *BGP<sub>i</sub>* with the *lowest cost* from remaining *BGPs* in *Q*Then search the *linked BGPs* to *BGP<sub>i</sub>*:

Linked BGPs are sorted by

- Number of shared variables
- Number of linked BGPs

• Filters and Values added to linked BGPs

### RESULTS

45000 r			5000 <sub>r</sub>		
40000	Hybrid+Sorting	Demographic and	4500		
40000	FedX FedX	geographic data (INSEE)			
35000		P1 (duplication)	4000		
		P2 (global rewriting)	3500		,
				1	







