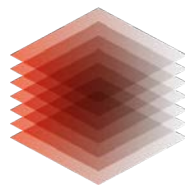




Leibniz
Universität
Hannover



TIB

SHACL Constraint Validation during SPARQL Query Processing

Philipp D. Rohde

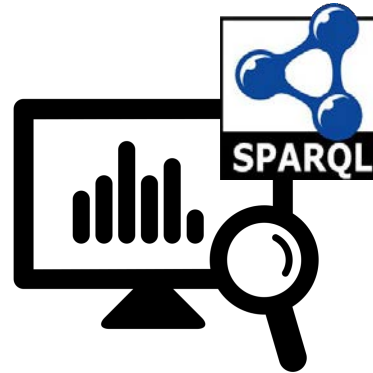
VLDB 2021 PhD Workshop

August 16th, 2021

Introduction



Knowledge Graphs (KGs)
gain Momentum



Data Retrieval:
SPARQL Queries



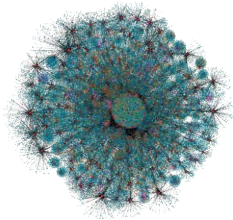
Integrity Constraints (ICs)
on RDF KGs

SHACL is the standard to specify ICs on RDF KGs

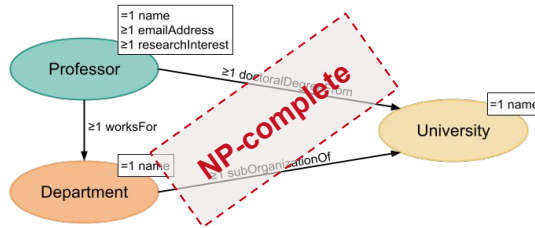
Motivating Example

INPUT

KG of a University System with 37,419 entities (~1M triples)



Shape Schema: Integrity Constraints on the KG



SPARQL Query: Retrieving Data from the KG

```
SELECT ?name ?ri ?uni WHERE {
  ?prof rdf:type ub:FullProfessor ;
  ub:name ?name ;
  ub:worksFor ?dept ;
  ub:doctoralDegreeFrom ?uni ;
  ub:emailAddress ?email ;
  ub:researchInterest ?ri .
}
```

- find physical plan
- execute SHACL validation
- produce results efficiently
- new optimization techniques based on the context are needed

OUTPUT

SPARQL Query Result with SHACL Validation Result Annotation

name	ri	uni	__meta__
FullProfessor0	Research6	http://www.Univeristy6.edu	all requirements met
FullProfessor3	Research10	http://www.University888.edu	University888 violates name constraint
...			

Query result annotation requires

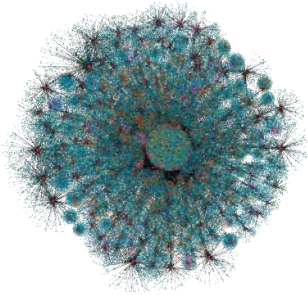
- SPARQL query execution
- SHACL shape schema validation

and is **computationally complex**

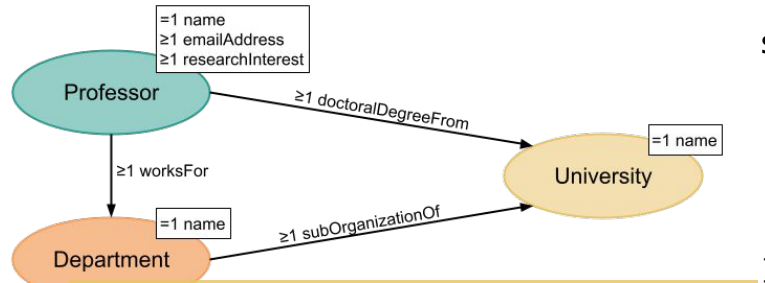
Motivating Example

INPUT

KG of a University System with 37,419 entities (~1M triples)



Shape Schema: Integrity Constraints on the KG



SPARQL Query: Retrieving Data from the KG

```

SELECT ?name ?ri ?uni WHERE {
  ?prof rdf:type ub:FullProfessor ;
  ub:name ?name ;
  ub:worksFor :Dept0 ;
  ub:doctoralDegreeFrom ?uni ;
  ub:emailAddress ?email ;
  ub:researchInterest ?ri .
}
  
```

Query result annotation requires

- SPARQL query execution
- SHACL shape schema validation

and is computationally complex

OUTPUT

name	ri	uni	__meta__
FullProfessor0	Research6	http://www.Univeristy6.edu	all requirements met
FullProfessor3	Research10	http://www.University888.edu	University888 violates name constraint
...			

Presentation Outline

1. Related Work
2. Proposed Approach
 - 2.1. Problem Statement
 - 2.2. Proposed Solution
 - 2.3. Online vs Offline Validation
3. Preliminary Results
 - 3.1. Trav-SHACL
 - 3.2. SPARQL Query Result Annotation
4. Lessons Learned
5. Conclusions

Related Work

SHACL Validation

- Complexity analysis of SHACL
 - NP-complete
- Identification of tractable fragments
 - no negation, but disjunction
 - no negation in recursion
 - no recursion
- SHACL validation using SPARQL endpoints
 - before: in-memory knowledge graphs
 - Datalog-like rules
- Recursive SHACL
 - left open in the specification

Integrity Constraints in Query Processing

Abbas et al.

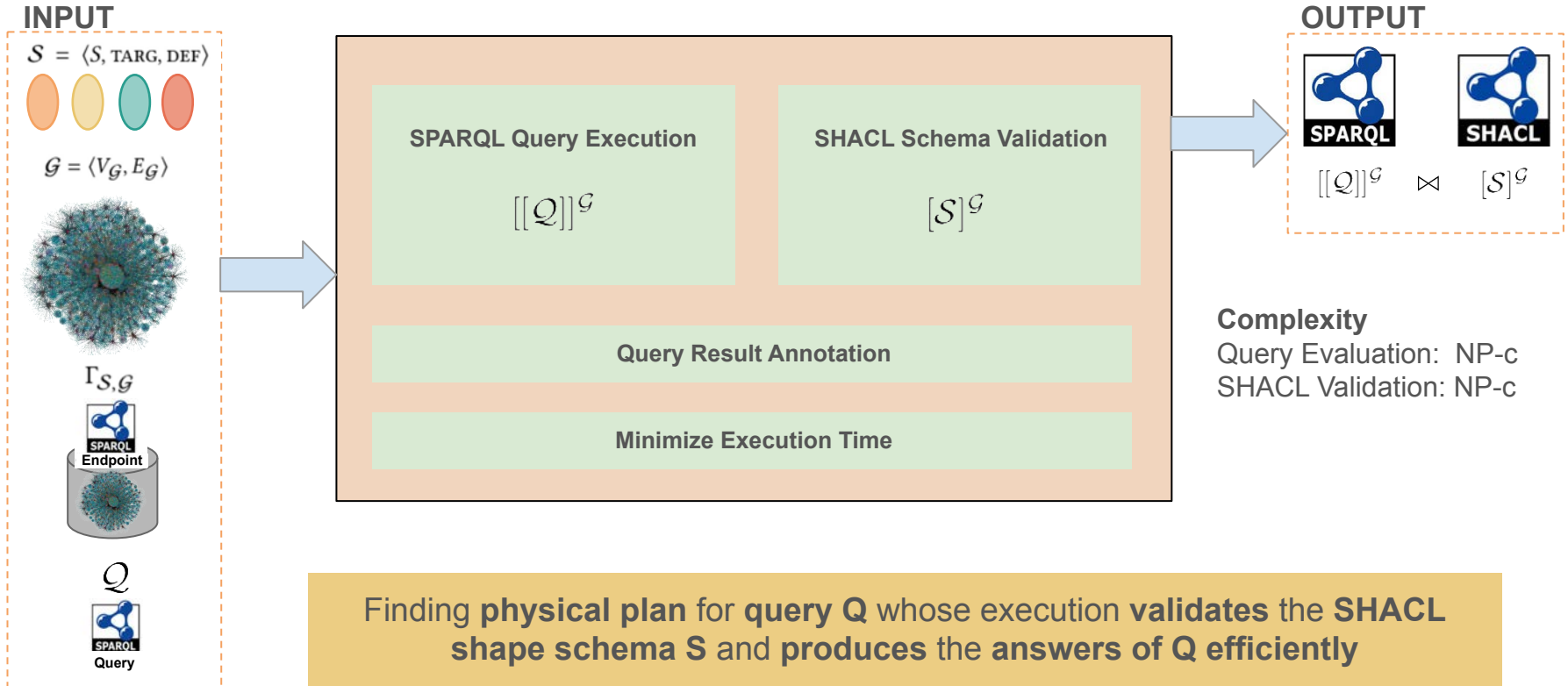
- well-formed ShEx schemas
- SPARQL triple pattern reordering
 - hierarchical structure of ShEx shapes
 - shapes included in other shapes are ranked higher
 - triple patterns with unique predicates are

No work on explainable SPARQL query results so far.

Rabbani et al.

- Extension of SHACL with statistics
- Cost-based query optimizer
- Precomputation time reduced

Problem Statement



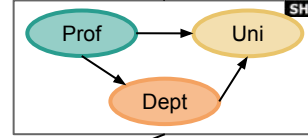
Proposed Approach

Query Decomposition

- subject star-shaped decomposition
- one star \approx one class

```

SELECT ?name ?ri ?uni WHERE {
  ?prof rdf:type ub:FullProfessor ;
  ub:name ?name ;
  ub:worksFor :Dept0 ;
  ub:doctoralDegreeFrom ?uni ;
  ub:emailAddress ?email ;
  ub:researchInterest ?ri .
}
    
```



SHACL Validation

- interleaved validation
- subset of shape schema

Query Result Annotation

- add SHACL validation result as metadata
- explainability

Novelty of the approach:

- identification of query plan able to combine query answering with integrity constraint validation
- explainability of SPARQL query results
- optimizations in SHACL validation

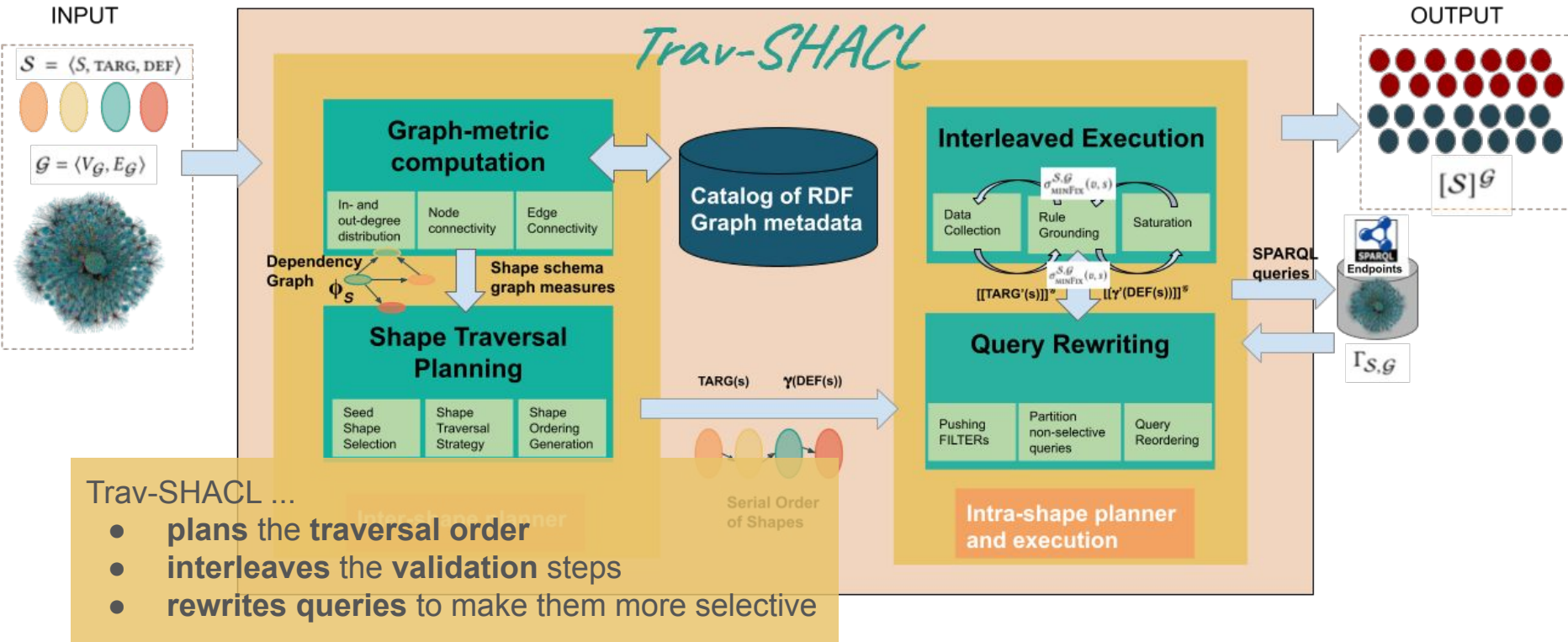
FullProfessor3 Research10 http://www.University888.edu University888 violates name constraint

Online vs Offline Validation

	online	offline
speed	✗	✓
who	✓ everyone	✗ data provider
adaptivity	✓	✗

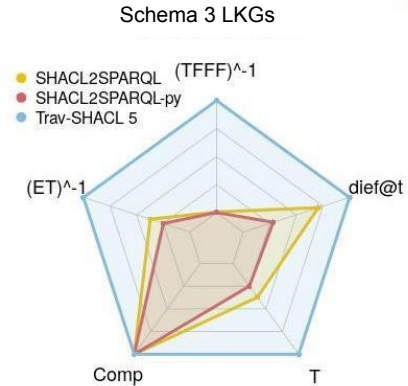
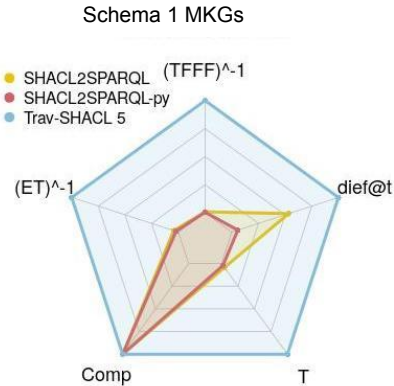
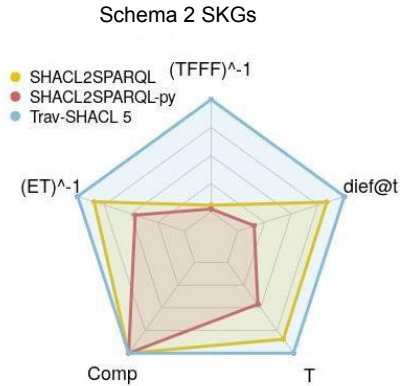
Offline validation is no option in this scenario

Results So Far: Trav-SHACL (1/2)



- Trav-SHACL ...
- plans the traversal order
 - interleaves the validation steps
 - rewrites queries to make them more selective

Results So Far: Trav-SHACL (2/2)



Metrics

dief@t: continuous efficiency at time t

(TFFF)⁻¹: Time for First Answer (sec)

(ET)⁻¹: Execution Time (sec)

Comp: sum of (in)validated entities

T: Throughput (answer/sec)

- # Constraint query mappings:
 - **839K** in SHACL2SPARQL,
 - **468K** in Trav-SHACL.

- # Constraint query mappings:
 - **703K** in SHACL2SPARQL,
 - **814** in Trav-SHACL.

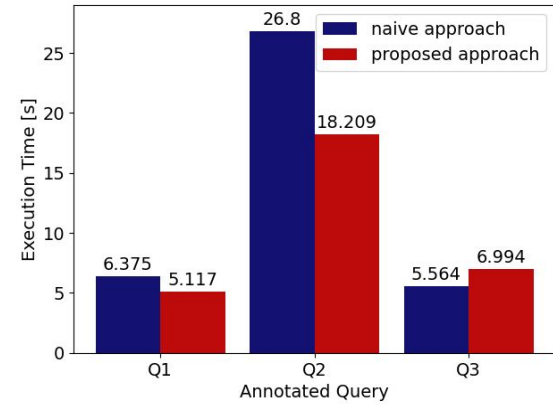
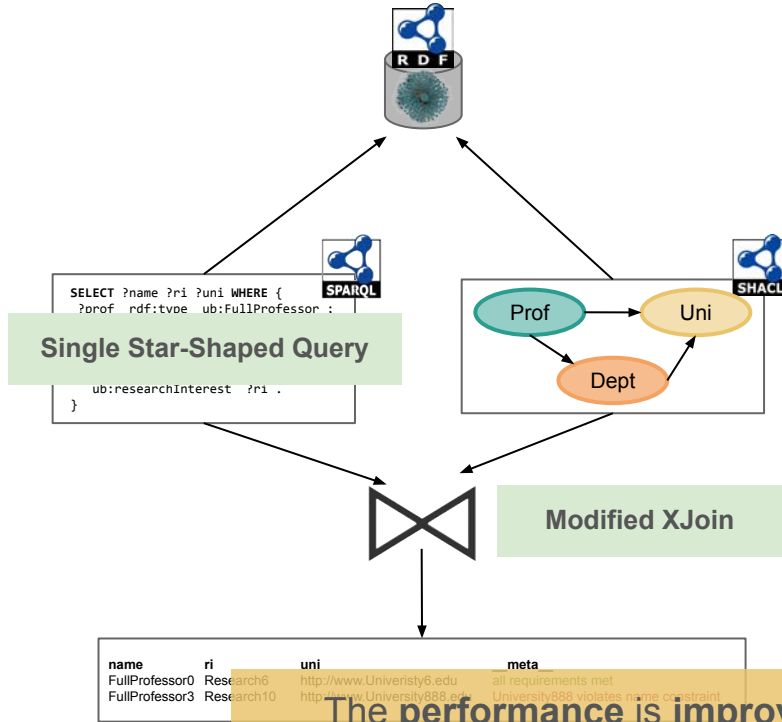
- # Constraint query mappings:
 - **22.94M** in SHACL2SPARQL,
 - **6.19M** in Trav-SHACL.

- Trav-SHACL always delivers results continuously,
- generates the first answer faster,
- finishes the execution faster,
- scales up to large knowledge graphs.



Impact of the interleaved execution

Results So Far: Query Annotation



RQ: Is the performance improved by applying the proposed approach?

WatDiv

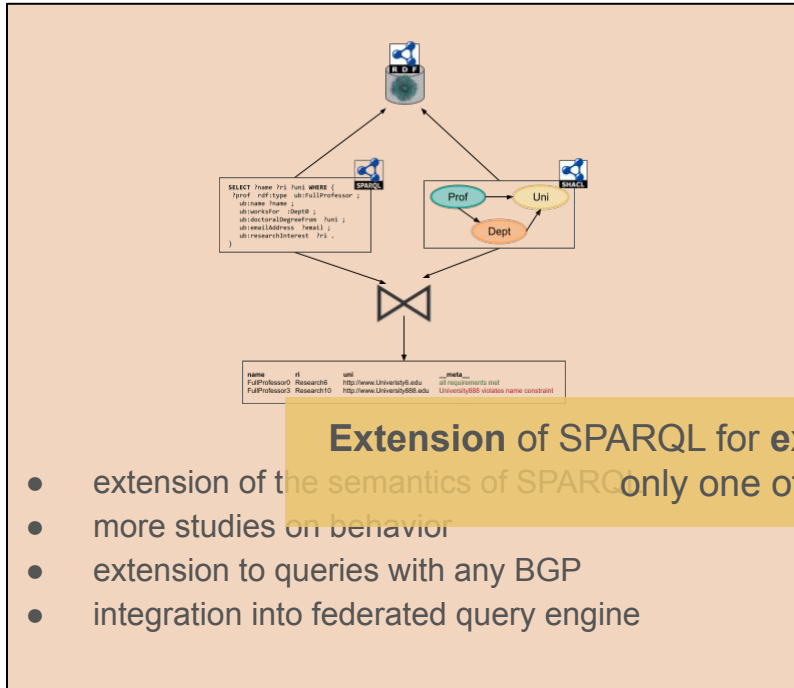
- 10 million triples
- 3-5 triple patterns per query
- less than 100 query results



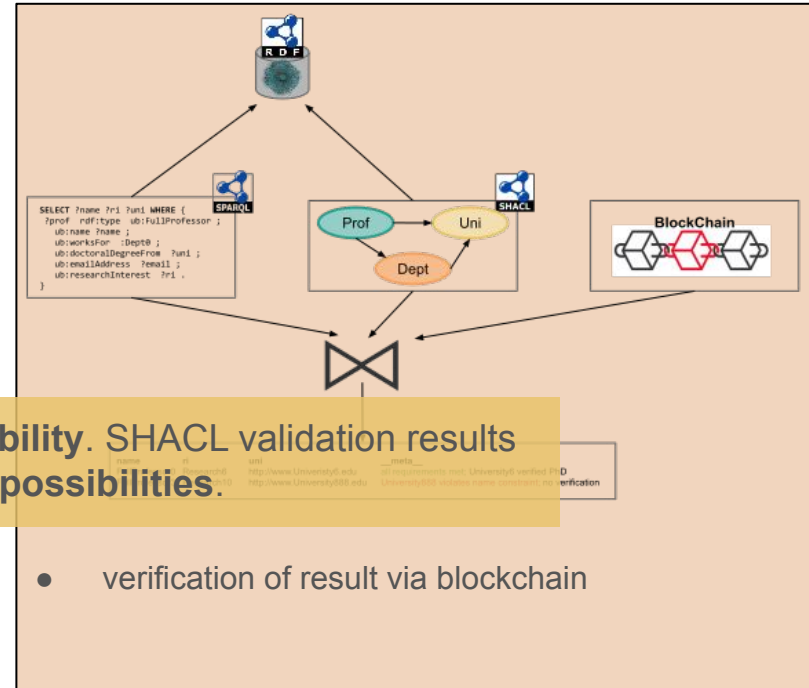
The performance is improved but more studies are needed.

Research Plan

SPARQL Query Result Annotation



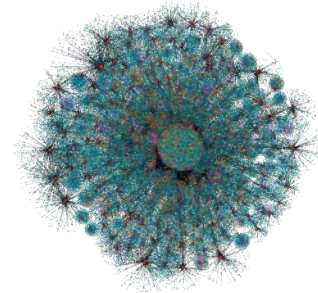
Additional Annotations



- extension of the semantics of SPARQL
- more studies on behavior
- extension to queries with any BGP
- integration into federated query engine

- verification of result via blockchain

Lessons Learned



- Benchmark needed
- Validation computationally expensive

- Explainability needed
- Lower performance is the price

- Many KGs have low data quality
- Means to improve quality are needed

Conclusions

Motivating Example

INPUT
 KG of a University System with 37,419 entities (~1M triples)

Shape Schema: Integrity Constraints on the KG

```

  class Prof {
    rdfs:type rdfs:Class
    rdfs:subClassOf Department
  }
  class Dept {
    rdfs:type rdfs:Class
  }
  class Uni {
    rdfs:type rdfs:Class
  }
  class FullProfessor {
    rdfs:type rdfs:Class
  }
  class DoctoralDegree {
    rdfs:type rdfs:Class
  }
  class ResearchInterest {
    rdfs:type rdfs:Class
  }
  Prof rdfs:subClassOf FullProfessor
  Prof rdfs:subClassOf DoctoralDegree
  Prof rdfs:subClassOf ResearchInterest
  Prof rdfs:subClassOf Dept
  Prof rdfs:subClassOf Uni
  
```

SPARQL Query: Retrieving Data from the KG

```

  SELECT ?name ?ri ?uni WHERE {
    ?prof rdf:type ub:FullProfessor ;
    ub:name ?name ;
    ub:worksfor ?dept ;
    ub:doctoralDegreeFrom ?uni ;
    ub:emailAddress ?email ;
    ub:researchInterest ?ri .
  }
  
```

OUTPUT

SPARQL Query Result with SHACL Validation Result Annotation

name	ri	uni	meta
FullProfessor0	Research6	http://www.University6.edu	all requirements met
FullProfessor3	Research10	http://www.University688.edu	University688 violates name constraint
...			

Page 2

Explainable SPARQL Query Results

Proposed Approach

Query Decomposition

- subject star-shaped decomposition
- one star = one class

SHACL Validation

- interleaved validation
- subset of shape schema

Query Result Annotation

- add SHACL validation result as metadata
- explainability

Page 4

Annotation with SHACL Validation

Results So Far: Query Annotation

Single Star-Shaped Query

Modified XJoin

Performance Comparison

Query	Native Approach (Execution Time)	Proposed Approach (Execution Time)
Q1	8.335	3.113
Q2	20.6	13.206
Q3	6.368	4.384

RQ: Is the performance improved by applying the proposed approach?

WatDiv

- 10 million triples
- 3-5 triple patterns per query
- less than 100 query results

Dataset Statistics

Category	Constraints	Instances
Dept	32 constraints	100,000 instances
Product	24 constraints	25,000 instances
Genre	2 constraints	145 instances

Page 11

Promising Results from Prototype

Research Plan

SPARQL Query Result Annotation

Additional Annotations

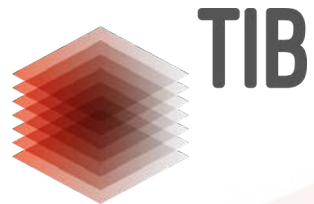
- more studies on behavior
- extension to queries with any BGP
- formalization of a SPARQL extension
- integration into federated query engine
- verification of result via blockchain

Page 10

Extension of SPARQL for Explainability



Leibniz
Universität
Hannover



Thanks for your attention!

Contact:

Philipp D. Rohde



philipp.rohde@tib.eu



[@philipp_rohde](https://twitter.com/philipp_rohde)