

DEPARTMENT OF

# **SIRIUS** Optique **SCIENCE**

### Faceted Search over RDF and OWL2

#### **Evgeny Kharlamov**

**Senior Research Fellow** Information Systems Group **Department of Computer Science** University of Oxford

Join work with

Marcelo Arenas Bernardo Cuenca Grau Sarunas Marciuska **Evgeny Sherkhonov** Dmitry Zheleznyakov

SemFacet Q

#### 2

#### Examples

JNIVERSITY OF

Media: BBC

enterprises

research labs

Geographic: Ordnance Survey

governmental institutions

- Governmental: Statistics.Data.gov.uk
- Industry: Siemens, Statoil, Aibel
- Knowledge Graphs: Google, DBpedia

### Semantic Data and Ontologies

RDF data and OWL ontologies are widely available in









SIEMENS





SemFacet @



### How to Access Semantic Data?

#### **SPARQL** Queries

- Not for End Users
- Like programming language: takes time to learn
- Error prone
- Complex error messages



SELECT	?z
WHERE	{?x rdf:type "president"@en;
	?x has_child ?z;
	{{?z is_graduated_from wiki:Harvard;}
	UNION
	<pre>{?z is_graduated_from wiki:Stanford;}</pre>

```
Virtuoso 37000 Error SP030:
SPARQL compiler, line 4: syntax error
at '.' before '}'
```



## Many Tools are Available to Access Sem Data

#### **Query formulation tools**

- Expressiveness
  - complex SPARQL queries
     ←→ exploration
- Control
  - precise
     ←→ ambiguous
- Usability
  - IT specialists
     ←→ end-user
- Adaption
  - Google
     ←→
     research prototypes
- Domain of use
  - Oil&Gas ← → Web





### **Faceted Search**

# Query formulation paradigm for naïve users

#### A facet = query constructor

- Name
- Set of values

#### **Facets in action**

- Choose a value
- Restrict search result





### Faceted Search: De Facto Web Query Standard





#### SemFacet

UNIVERSITY OF

FORI

### **Classical Faceted Search: Data Model**

#### Data

- Set of docs
- Annotated with strings

#### Example

- documents:
  - politician
- annotations:
  - type:
    - President,
       Revolutionary
  - citizen of
    - RUS, US
  - has child

- Yes, No







### **Classical Faceted Search: Search Process**





and George H. W. Bush, he was born in New Haven,

Connecticut. After graduating from Yale University in ...

acet @

## Faceted Search over RDF: Data Model

#### Data

- Graph
- Documents (URIs)
  - of different kind
  - annotated with
    - strings
    - other URIs

#### Example

- documents:
  - Politician
  - Countries
  - Children
  - Universities
- annotations:
  - Type: President, Revolutionary





### Faceted Search over RDF: Search Process



## **Existing Semantic Faceted Search Solutions**

#### **Systems**

- Parallax
- /facet
- BroweRDF
- F-Search in Information Workbench
- mSpace
- Tabulator
- Humboldt
- GoPubMed
- gFacet
- ...

UNIVERSITY OF





## Were We Happy Existing Solutions?

#### Existing solutions (we aware of) are RDF data driven

- they (essentially) ignore ontological axioms
- they ignore complex (OWL 2) data statements

#### **Theoretical underpinnings**

- have received very little attention
- What fragments of SPARQL can be naturally captured?
- What is the complexity of answering such queries?
  - when data is enhanced with ontologies in OWL 2 profiles
- How can we formally capture interactive interface update?



D

UNIVERSITY OF

Interface

12



### **Our Contributions**

#### Formalise faceted interfaces tailored towards RDF and OWL

- Capture key functionalities in implemented systems
- Abstract from GUI-specific considerations

#### Study of expressive power and complexity of faceted queries

- Identify a fragment of FO sufficient to capture such queries
- Study complexity of query evaluation wrt ontologies in the OWL 2 profiles

#### Study interface generation and update

- Lift existing approaches RDF  $\rightarrow$  OWL 2
- Propose generic algorithms

#### **Develop SemFacet system**

- Allows to generate faceted search interfaces over RDF and OWL 2
- Scales to millions of triples





### Facet: Formally

#### What is Facet?

- an elementary building block of faceted interfaces
- it says what can be selected

#### Facet consists of three components

- Property
- Logical OR or AND
- Set of values,
  - e.g., classes, objects, strings





#### Facets:

- $F_1 = (\mathsf{type}, \lor \{\mathsf{Pres}, \mathsf{Revol}\})$
- $F_2 = (\mathsf{haschild}, \lor \{\mathsf{any}, d_{kr}, d_{cc}, d_{bb}, d_{ep}\})$
- $F_3 = (\mathsf{gradFrom}, \vee \{d_{hu}, d_{su}, d_{yu}, d_{spu}\})$
- $F_4 = (\mathsf{citizenOf}, \lor \{d_{ru}, d_{us}\})$

#### http://en.wikipedia.org/wiki/Bill Clinton

William Jefferson "Bill" Clinton (born William Jefferson Blythe III; August 19, 1946) is an American politician who served as the 42nd President of the United States from 1993 to 2001. Inaugurated at age 46, he was the third-youngest president. He took office at the end of the Cold War, and was the first president of the baby boomer generation...

#### http://en.wikipedia.org/wiki/Theodore\_Roosevelt

Theodore "Teddy" Roosevelt was the 26th President of the United States (1901–1909). He is noted for his exuberant personality, range of interests and achievements, and his leadership of the Progressive Movement, as well as his "cowboy" persona and robust masculinity...

#### SemFacet 🕲

Search

### Faceted interface: Formally

#### **Faceted Interface**

- Basic Faceted Interface
  - Facet + Set of selected values
- General Faceted Interface:
  - Bool, combination of nested BIs

#### According to a grammar $I ::= path \mid (path \land path) \mid (path \lor path),$ path ::= $I_0 | (I_1/I)$ .

### **Nested Faceted Interface:** $(F_1, \{\mathsf{Pres}\})$ $(F_2, \{any\})/(F_3, \{d_{su}, d_{hu}\})$ $\wedge$ $(F_4, \{d_{us}\})$



#### Facets:

- $F_1 = (type, \lor \{Pres, Revol\})$
- $F_2 = (\mathsf{haschild}, \lor \{\mathsf{any}, d_{kr}, d_{cc}, d_{bb}, d_{ep}\})$

$$F_3 = (\mathsf{gradFrom}, \lor \{d_{hu}, d_{su}, d_{yu}, d_{spu}\})$$

$$F_4 = (\mathsf{citizenOf}, \lor \{d_{ru}, d_{us}\})$$

http://en.wikipedia.org/wiki/Bill Clinton William Jefferson "Bill" Clinton (born William Jefferson Blythe III; August 19, 1946) is an American politician who served as the 42nd President of the United States from 1993 to 2001. Inaugurated at age 46, he was the third-youngest president. He took office at the end of the Cold War, and was the first president of the baby boomer generation...

#### http://en.wikipedia.org/wiki/Theodore\_Roosevelt Theodore "Teddy" Roosevelt was the 26th President of the United States (1901–1909). He is noted for his exuberant personality, range of interests and achievements, and his leadership of the Progressive Movement, as well as his "cowboy" persona and robust masculinity...



### SemFacet @

Search

### Faceted interface with Refocusing

#### How we formalized refocusing

 Essentially: special value with special treatment

#### Facet with focus:

 $F_2 = (\mathsf{haschild}, \lor \{\mathsf{any}, d_{kr}, d_{cc}, d_{bb}, d_{ep}\})$  $F_5 = (\mathsf{haschild}, \lor \{\mathsf{focus}, d_{kr}, d_{cc}, d_{bb}, d_{ep}\})$ 

#### **Nested Faceted Interface:**

 $(F_1, \{\mathsf{Pres}\})$   $\land$  $(F_5, \{\mathsf{focus}\})/(F_3, \{d_{su}, d_{hu}\})$   $\land$  $(F_4, \{d_{us}\})$ 





### **Turning Nested Faceted Interfaces into Queries**



### **Faceted Queries**

#### Faceted queries are FO formulas of restricted shape

- Positive existential formulas
- Monadic
- Variables form a directed tree rooted at the free variable
- Every two formulae disjunctively connected share one variable

Pres x  $\bigcirc \bigcirc /(hasChild(x, y))$   $\bigcirc /(hasChild(x, y))$  $\bigcirc$ 



Dependency Tree



### Faceted query answering

#### **Over RDF datasets**

- Tractable
- Tree shaped queries with restricted disjunction
- Bottom up evaluation
- as opposed to NP-Hard for unrestricted positive existential queries



### Interface Generation & Update w/ F-graph



## Navigation Graph



#### Facet graph: our main "data structure" for interface gen. & update:

- Nodes are possible facet values (unary predicates and constants)
- Edges are possible facet names (binary predicates or type)
- Every edge must be justified by
  - An entailed fact
  - an entailed axiom
    - projection of axioms on a graph





### **Interface Generation & Update**



#### Our generation and update algorithms are

- "guided" by
  - data
  - ontology
- justified by entailments
  - each facet & value in initial interface
  - each interface update is justified by entailments



### **Interface Generation & Update**



- justified by entailments
  - each facet & value in initial interface
  - each interface update is justified by entailments



### SemFacet System

#### Integration of

- Keyword search and
- faceted search over RDF & OWL 2

#### Main features

- Automatic generation of faceted interfaces
- In memory
- Online and offline reasoning
- Efficient on millions of triples
- Backend: RDFOX, PAGOdA, Hermit, Sesame

#### Flexible configuration

- Interchangeable triple stores
- Configurable answers (snippets)
- Support of Or and And facets





Statistics		Active Settings	
Active settings	Max results returned by the keyword search:	1000	
Data Preview	Snippet image predicate:	thumbnail	
	Snippet url predicate:	isPrimaryTopicOf	
Jpdate settings	Snippet description predicate:	comment	
oad data	Snippet description separator:	£	
SemFacet Search	Facet category predicate:	rdf.type	
	Facet category relation predicate:	rdfs:subClassOf	
	Facet nesting:	disabled	
	Conjunctive facet predicates:	hasRating, isLocatedIn	

## Contributions

#### Foundations

- Formalization of Semantic Faceted Search
- Complexity of query answering

#### Projection of ontologies on graphs

- Allows to incorporate ontologies into faceted search
- Better faceted interfaces: generate more facets / prune irrelevant facets

#### Scalable algorithms to

- Generate and update facets from data and ontologies
- Evaluate faceted queries over semantic data

#### System implementation

SemFacet System





#