

SQUALL: a Controlled Natural Language for Querying and Updating RDF Graphs

Sébastien Ferré
Team LIS, Data and Knowledge Management, Irisa

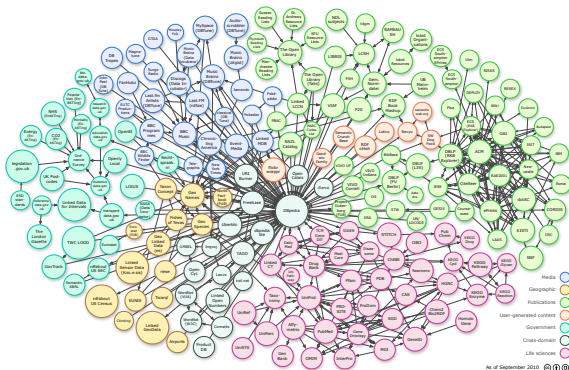
Controlled Natural Language, 30 August 2012, Zurich

INSTITUT DE RECHERCHE EN INFORMATIQUE ET SYSTEMES ALÉATOIRES



The Web of Data

- ▶ How to search and explore the **Web of data** (RDF graphs) ?
- ▶ How to fill the gap between **end users** and **formal languages** (RDF, OWL, SPARQL) ?



Formal vs Natural Languages

- ▶ SPARQL: a formal language *à la* SQL
 - ▶ very expressive and precise for querying and updating RDF graphs
 - ▶ requires understanding of low-level notions: relational algebra and logic
- ▶ natural language interfaces (ex., Aqualog, FREyA)
 - ▶ good usability through NL
 - ▶ difficult problems: ambiguity and adequacy w.r.t. the underlying system
 - ▶ in practice, generally limited to simple questions (much less expressive than SPARQL)
 - ▶ ex., Aqualog queries are limited to 2-triples queries



Controlled Natural Languages (CNL)

- ▶ on the natural/formal continuum [Kaufmann&Bernstein 2010]
- ▶ combine **natural syntax** and **formal semantics**
- ▶ *“There is no important theoretical difference between natural languages and the artificial languages of logicians.”* (Montague)
- ▶ a few CNLs:
 - ▶ **ACE** [Fuchs *et al*]: a general purpose CNL
 - ▶ **SOS**, **Rabbit**: CNLs for verbalizing OWL axioms
 - ▶ **SQUALL**: the first CNL for SPARQL queries *and* updates



What SQUALL is *not*...

1. a pure (grammatically correct) **subset** of English
 - ▶ natural languages are a source of inspiration for flexibility, expressiveness, concision, high-level forms
 - ▶ I think that CNLs should be **more regular** than NLs because they have to be learnt anyway
2. concerned with **morphology** (lexicon, agreements, etc.)
 - ▶ should have the same requirements as SPARQL w.r.t. data
 - ▶ should be able to refer to every resource without preprocess
 - ▶ shares **non-ambiguous notations** with SPARQL
 - ▶ `author` → `hackcraft:authoredBy` **or** `purl:author` ... ?
3. a **user interface**
 - ▶ querying is difficult, whatever the language
 - ▶ syntax errors, empty results, preferences
 - ▶ syntactic guided input (e.g., Ginseng) is not enough
 - ▶ the objective is **semantic guided input** (done to a limited extent in previous work with Sewelis)



What SQUALL is...?

- ▶ an alternative CNL syntax for SPARQL
 - ▶ hence, same expressiveness as SPARQL
 - ▶ hence, full adequacy to RDF data
 - ▶ with natural high-level syntax
- ▶ its implementation is a compiler (3 phases)
 1. parsing of the source sentence (SQUALL)
 2. generation of an intermediate representation (Montague λ -terms)
 3. production of the target code (SPARQL)



RDF Graphs

- ▶ a **RDF graph** is a set of triples (labeled edges)
- ▶ a **triple** (of resources) has a subject, a predicate, and an object
 - ▶ a triple is a basic sentence
 - ▶ **ex.**, `ex:John ex:loves ex:Mary .`
- ▶ a **resource** denotes an entity or a concept or a literal value (e.g., numbers, dates, strings)
- ▶ a **property** denotes a binary relation between resources
 - ▶ a property can be a transitive verb (ex., `ex:loves`) or a relational noun (ex., `ex:author`)
- ▶ a **class** denotes a set of resources
 - ▶ a class can be a noun (ex., `ex:woman`) or an intransitive verb (ex., `ex:works`)
- ▶ properties and classes are resources themselves



SPARQL 1.1: querying and updating RDF graphs

- ▶ SPARQL forms
 - ▶ closed question: `ASK graph pattern`
 - ▶ open question: `SELECT vars WHERE graph pattern`
 - ▶ update: `DELETE graph INSERT graph WHERE graph pattern`
- ▶ graph patterns
 - ▶ relational algebra: joins, unions, complements, selections, projections
 - ▶ constraints: logic, arithmetic, built-ins
 - ▶ named graphs
 - ▶ subqueries
 - ▶ aggregations



A Montague grammar of SQUALL

- ▶ Syntax and semantics
- ▶ Modules:
 1. lexical conventions
 2. triples as sentences
 3. relational algebra as coordinations
 4. natural constructs (headed NPs, relatives, ...)
 5. queries with *wh*-words
 6. quantifiers as determiners
 7. subordination and n-ary predicates with reification
 8. built-in predicates and aggregations
 9. resolving syntactic ambiguities



Lexical conventions

The same as in well-known notations (Turtle, SPARQL, N3)

- ▶ proper nouns, nouns, and verbs (URIs)
 - ▶ `<http://dbpedia.org/resource/Berlin>`: a full URI for the Berlin city
 - ▶ `dbpedia:Berlin`: an abbreviated URI with DBpedia namespace
 - ▶ `:Berlin`: an abbreviated URI with default namespace
 - ▶ `Berlin`: a bare URI (default namespace)
- ▶ literals
 - ▶ `"Hello world!"`: a plain literal
 - ▶ `"42"^^xsd:integer`: a typed literal
 - ▶ `42`: a bare integer
- ▶ variables: `?X`
- ▶ grammatical words (SQUALL reserved keywords)
 - ▶ `is`, `a`, `which`, `every`, ...



Relational algebra as coordinations

They apply to most syntagms Δ : *S*, *NP*, *VP*, *P1*, *P2*, (next: *Rel*, *AP*, *PP*).

Δ	\rightarrow	not Δ_1	{	not δ_1	}		“not [<i>VP</i> know-s B]”	
		Δ_1	and	Δ_2	{	and δ_1 δ_2	}	“[<i>VP</i> work-s] and [<i>VP</i> cite-s X]”
		Δ_1	or	Δ_2	{	or δ_1 δ_2	}	“[<i>NP</i> A] or [<i>NP</i> B]”
		maybe Δ_1	{	option δ_1	}		“maybe [<i>VP</i> know-s B]”	



Built-in predicates and aggregations

- $P1 \rightarrow Pred1URI \{ \lambda x.(\mathbf{pred1} \text{ uri } x) \}$ “Monday”
- $P2 \rightarrow Pred2URI \{ \lambda x.\lambda y.(\mathbf{pred2} \text{ uri } x \ y) \}$ “match”
- $NG1 \rightarrow \text{AggregURI of AP (per } AP_i^+ \text{)?}$
 $\{ \lambda x.(\mathbf{aggreg} \text{ uri } x \text{ ap } (ap_i)_i) \}$
“count of [AP the publication ?P] per [AP the year of ?P]”



Translation to SPARQL (principle)

Starting with the λ -term produced when parsing a SQUALL sentence:

1. replace some constants by their definition
 - ▶ ex., **count** = $\lambda d.(\mathbf{select} \lambda x.(\mathbf{aggreg} \text{COUNT } x \ d \ ()))$
2. perform all β -reductions on the result
3. inductively generate SPARQL code
 - ▶ 4 generators for: sentences, updates, queries, and graph patterns
 - ▶ $[\mathbf{ask} \ f] = \text{ASK} \{ [f]_G \}$
 - ▶ $[\mathbf{select} \ d] = [?x \mid d \ ?x]_Q$
 - ▶ $[\mathbf{forall} \ d_1 \ d_2]_G = [\mathbf{not} \ (\mathbf{exists} \ (\mathbf{and} \ d_1 \ (\mathbf{not} \ d_2)))]_G$
 - ▶ $[\mathbf{the} \ d_1 \ d_2]_G = [\mathbf{exists} \ (\mathbf{and} \ d_1 \ d_2)]_G$
 - ▶ $[\mathbf{the} \ d_1 \ d_2]_U = [\mathbf{forall} \ d_1 \ d_2]_U$



Translation to SPARQL (example)

The SQUALL sentence

for which researcher-s ?X, in graph DBLP every publication whose author is ?X and whose year \geq 2000 has at least 2 author-s

is parsed as

“*[Sfor [NP[Detwhich] [NG1[P1researcher-s] [AR[App?X]]]],
[S[PPin [Prepgraph] [NPDBLP]] [S[NP[Detevery]
[NG1[P1publication] [AR[Rel[Relwhose [NG2[P2author]] [VPis
[AP?X]]] and [Relwhose [NG2[P2year]] [VP[P2 \geq] [NP2000]]]]]]]]*”



Translation to SPARQL (example)

...whose internal representation is:

```
(select X (and
  (triple X rdf:type :researcher)
  (graph :DBLP
    (forall (exists x3 x5 (and
      (triple x3 rdf:type :publication) (triple x3 :author X)
      (triple x3 :year x5) (pred2 ≥ x5 2000)))
    (exists x6 (and
      (aggreg COUNT x6 x8 x3
        (exists x8 (triple x3 :author x8)))
      (pred2 ≥ x6 2))))))))
```



Translation to SPARQL (example)

...which translates to the SPARQL query:

```
SELECT ?r
WHERE {
  ?r rdf:type :researcher .
  BIND (?r AS ?X)
  GRAPH :DBLP {
    FILTER NOT EXISTS {
      ?p rdf:type :publication .
      ?p :author ?X .
      ?p :year ?y .
      FILTER (?y >= 2000)
      FILTER NOT EXISTS {
        { SELECT COUNT(?a) AS ?n
          WHERE { ?p :author ?a . } }
        FILTER (?n >= 2) } } } }
```



Perspectives as a language

- ▶ covering 100% of SPARQL 1.1 (nothing hard)
 - ▶ CONSTRUCT and DESCRIBE forms
 - ▶ query modifiers (ORDER BY, LIMIT, OFFSET)
 - ▶ conditional expressions
 - ▶ concise notation of RDF collections (list) and membership test
 - ▶ + type checking in expressions
- ▶ adding more natural constructs
 - ▶ anaphoras: e.g., “some man ... this man” instead of “some man ?X ... ?X”
 - ▶ comparatives and superlatives
 - ▶ adjectives and adverbs (as RDF classes)
 - ▶ more forms of aggregations
 - ▶ e.g., “what is the average age of the researcher-s ?”



Perspectives as a user interface

- ▶ even a CNL is not simple enough
- ▶ need for guided input
 - ▶ syntax-based auto-completion (like in Ginseng)
 - ▶ query-based faceted search (like in Sewelis)



