

Who the FOAF knows Alice? A needed step toward Semantic Web Pipes ^{*}

Giovanni Tummarello¹, Axel Polleres¹, and Christian Morbidoni²

¹ DERI Galway, National University of Ireland, Galway
{firstname.lastname}@deri.org

² SeMedia Group, Universita' Politecnica delle Marche, Ancona, Italy
christian@deit.univpm.it

Abstract. In this paper we take a view from the bottom to RDF(S) reasoning. We discuss some issues and requirements on reasoning towards effectively building Semantic Web Pipes, aggregating RDF data from various distributed sources. If we leave out complex description logics reasoning and restrict ourselves to the RDF world, it turns out that some problems, in particular how to deal with contradictive RDF statements, do not yet find their proper solutions within the current Semantic Web Stack. Besides theoretical solutions which involve full DL reasoning, we believe that more practical and probably more scalable solutions are conceivable one of which we discuss in this paper, namely, expressing and resolving conflicting RDF statements by means of a specialized RDF merge procedure. We implemented this conflict-resolving merge procedure in the DBin system.

1 Introduction

Publishing RDF files on the Web is bound to become more and more a way to state facts that are asserted or believed to be true by the producer of the source itself. DBpedia [1], for example, publishes a large collection of such facts by extracting them from the collective works of the Wikipedia communities. FOAF [5] files are personal RDF models which are created by individuals to state facts about, typically, themselves. Nothing, however, prevents them in general to state facts about other entities and this is in fact a fundamental feature of the "Semantic Web", everyone is allowed to "state" about, virtually, anything. In some cases one might even be inclined to trust third-party information more than self-descriptions, for instance comments about an enterprise or a product one considers to buy. The sum of RDF statements, currently known to be HTTP retrievable, is now in the order of billions with millions of individual HTTP locations (sources) hosted on tens of thousands of web sites, rapidly increasing.

^{*} This work has been partially supported by the European FP6 project in-Context (IST-034718), by Science Foundation Ireland under the Lion project (SFI/02/CE1/I131), and by the European project DISCOVERY (ECP-2005-CULT-038206).

Along with this increased take-up of RDF on the Web, upcoming query language standards like SPARQL [12], or RDF search engines like SWSE [6] or Sindice [13] shall finally enable structured querying over Web data. Unfortunately however, there is no clear and established model on how to use such amounts of information coming from many diverse sources. Using any available source directly, e.g. crawling/downloading and using it might not be advisable or sufficient. More information might be needed such as, for example, patches to the original data. Other cases include when a source is in general considered useful but is known to contain statements which need to be removed, e.g. outdated facts (a "negative" information patch is needed), or subjective assertions which can be accepted or not depending on who is reading the data. In general, getting information from the Web into one's own semantic client or system is very likely to require, or at least benefit, from a series of custom steps to be performed involving a number of external or internal sources before having a version which can be used directly. Also, facing the sheer amount of data to be expected, more complex tasks such as ontological inferences or complex query answering will profit from such preprocessing which only preserves relevant and useful information. In this paper, we focus on one facet of such preprocessing, namely allowing to retract unwanted RDF data and present a practical solution for this problem.

2 Towards Semantic Web Pipes

Yahoo Web Pipes³ are a recent development which has certainly had already a big impact to the latest wave of web development by showing how customized services and information streams can be implemented by sequentially processing and interleaving existing feeds and services. With Pipes, resources, e.g. RSS feeds, can be merged with one another, filtered according to specific pipe rules, used as an input for an on-line restful API to get yet more results, etc. Most interestingly, this all happens without the original providers of informations and services had to change anything on their side or reach any form of agreement if not to use HTTP and possibly RSS. Current mashup models like Yahoo Pipes are however limited to "streams" of information (e.g. news feeds) or single, simple API invocations on a remote site (e.g. a search for a specific word, or, more general, one-shot Web service invocations).

In the same way as a Web Pipe enables an existing Web information stream to be customized, extended and reused for a specific purpose as decided by the pipe creator, we see a very clear interest in trying to use this model to address the issue we highlighted before: how to make use of web published RDF sources? We might for example want to use DBpedia knowledge about a topic, but yet sum it with the knowledge coming from certain specific sites and correcting it by eliminating some statements we believe to be false. The Web Pipe model teaches us that we do not really want to download the DBpedia RDF dump, and operate directly on a local version of it, e.g. by adding and subtracting triples in a complex SPARQL query (see also the following Section). By doing

³ <http://pipes.yahoo.com/>

so once and in a static manner, we would create a customized knowledge base at the beginning but would miss any new information that any of the composing sources might later add. A much more dynamic and useful model would therefore be a "Semantic Web Pipes" model where an RDF piping engine can on the fly and on demand work out the customized composition and processing of a set of Web sources according to our specific needs. In case where information needs to be simply added, the RDF semantics [7] specifies how to merge two models: the piping engine has therefore to do not much more than downloading the files and putting them together in the same store, standardizing apart blank nodes. But what to do when information needs to be patched in a traditional sense, i.e. in part both removed and added?

As a use case, let us take the case where Bob is stating that Charles knows Alice in his FOAF [5] file. Alice has a questionable reputation, and Charles, clearly, has no control on Bob's FOAF file. Clearly, a minimal requirement on distributed metadata is the ability to counter such false statements, thus giving Charles a way to state in his FOAF file a simple and unambiguous statement: "I don't know Alice". We aim to provide a simple and minimalistic solution to this problem, thus avoiding unnecessarily complex reasoning.

3 Related Works: Expressing Negative RDF Statements

First, we note that neither RDF nor RDF Schema provide means to make negative statements such as "Charles doesn't foaf:know Alice", see last statement in Figure 1(b).

<pre> @prefix : <http://examp.org/ bob#> @prefix foaf: <http://xmlns.com/foaf/0.1/> :me foaf:name "Bob". :me foaf:knows <http://alice.exa.org/i> . :me foaf:knows <http://ex.org/charles#me>. <http://ex.org/ charles#me> foaf:knows <http://alice.example.org/i>. ... </pre>	<pre> @prefix : <http://ex.org/ charles#> @prefix foaf: <http://xmlns.com/foaf/0.1/> @prefix rdf: <http://www...rdf-syntax-ns#> :me rdf:type foaf:Person; foaf:name "Charles". :me foaf:knows <http://examp.org/bob#me>. :me foaf:knows <http://alice.exa.org/i>. ... </pre>
(a) Bob's FOAF file	(b) Charles' FOAF file

Fig. 1. Personal information in FOAF

The semantics of RDF(S) is purely monotonic and described in terms of positive inference rules, so even if Charles added instead a new statement

```
:me myfoaf:doesntknow <http://alice.exa.org/i> .
```

he would not be able to state that statements with the property `myfoaf:doesntknow` should single out⁴ `foaf:knows` statements.

⁴ In fact, we mean here overriding instead of simply contradicting in the pure logical sense.

N3 Tim Berners-Lee's Notation 3 (N3) [2] provides to some extent means to express what we are looking for by the ability to declare falsehood over reified statements which would be written as:

```
{ :me foaf:knows <http://alice.exa.org/i> } a n3:falsehood .
```

Nonetheless, this solution is somewhat unsatisfactory, due to the lack of formal semantics for N3; N3's operational semantics is mainly defined in terms of its implementation *cwm*⁵ only.

OWL The falsehood of Charles knowing Alice can be expressed in OWL, however in a pretty contrived way, as follows (for the sake of brevity we use DL notation here, the reader might translate this to OWL syntax straightforwardly):

$$\{charles\} \in forall.foaf:knows \neg \{alice\}$$

Reasoning with such statements firstly involves OWL reasoning with nominals, which most DL reasoners are not particularly good at, and secondly does not buy us too much, as the simple merge of this DL statement with the information in Bob's FOAF file would just generate a contradiction, invalidating all, even the useful answers.

Para-consistent reasoning on top of OWL, such as for instance proposed in [8] and related approaches, solve this problem of classical inference, but still requiring full OWL DL reasoning.

SPARQL Finally, more along the Pipes idea, one could as a naive solution, deploy an off-the-shelf SPARQL engine and filter Bob's FOAF file by a query, leaving just the clean statements. Imagine that Charles stores his unwanted statements in the RDF Web source `<http://ex.org/~charles/badstatements.rdf>`, then such a query filtering the information from merging Bob's and Charles' FOAF files could look as follows:

```
CONSTRUCT { ?S ?P ?O }
FROM <http://ex.org/~charles/foaf.rdf>
FROM <http://ex.org/~bob/foaf.rdf>
WHERE { ?S ?P ?O .
        OPTIONAL { GRAPH <http://ex.org/~charles/badstatements.rdf>
                    { ?S1 ?P1 ?O1 . }
                  FILTER (?S1 = ?S && ?P1 = ?P && ?O1 = ?O &&) }
        FILTER ( !Bound(?S1) ) }
```

Simply putting the bad information in a separate file, is not a proper solution for the scenario we outlined, as it is not clear how a Crawler stumbling over `<http://ex.org/~charles/badstatements.rdf>` should disambiguate it from valid RDF information. Rather, we would need to reify the negative statements using for instance the N3 version outlined before, or the "native" RDF Reification vocabulary which would besides blowing up metadata by unhandy reified statements, further complicate the SPARQL query⁶ to filter out the "good" data.

⁵ <http://www.w3.org/2000/10/swap/doc/cwm>

⁶ Note that, in the FILTER query, we exploit the admittedly awkward way to model set difference in SPARQL which as such might already not be considered intuitive unanimously.

In the following, we will sketch a more practical solution to the problem, exploiting previous work on Minimum Self Contained Graphs.

4 Introducing "Negative Statements" using MSG Hashes

Any RDF graph may be viewed as set of triples. Triple level processing of distributed RDF files, particularly identifying the same RDF graphs, is made very complex by the existence of blank nodes. For this reason, the RDFSynC algorithm introduce the notion of Minimum Self Contained Graph (MSGs) [14]. Simply said, an MSG is constructed starting from a triple and collecting, for each blank node in it, all the other triples attached to these until no more blank nodes are involved. Such "closure" makes sure that a graph can be recomposed at a different location simply by merging all the MSGs by which it is composed, even if these are transferred one at a time. As MSGs are stand-alone RDF graphs, they can be processed with algorithms such as canonical serialization. We use an implementation of the algorithm described in [4], which is part of the RDFContextTools Java library⁷ to obtain a canonical string representing the MSG and then we hash it to an appropriate number of bits to reasonably avoid collisions. This hash acts as a unique identifier for the MSG with the fundamental property of being content based, which implies that two remote peers would derive the same ID for the same MSG in their DB. A graph can be therefore treated as a set of digital hashes each one representing an MSGs. In the context of the problem addressed in the present work, we use such digital hashes to refer to the MSG itself. This takes usually the form of a literal of encoding the 16 bytes of the MD5 hash. Stating that an MSG is false is therefore as easy as stating a single triple where the Subject is the source, the predicate is a designated predicate (we suggest for instance the URI `<http://sw.deri.org/09/2007/states_not>`) and the 16 bytes literal containing the MSG hash as an object, so the negative statement could be made directly in Charles' FOAF file or in a separate file as follows:

```
:me <http://sw.deri.org/09/2007/states_not>  
    "HASH OF :ME FOAF:KNOWS ALICE"^^msg:Hashsum .
```

Storing MSG hashes instead of reifying statements has (except saving storage space) some other interesting implications: This solution works conceptually well as it takes care even of the cases where one wants to deny statements which involve blank nodes. This would be possible using reification due to the arising ambiguity. Digital hashes over MSGs, which are agnostic about blank node IDs, avoid this problem.

One drawback of the solution to quasi "encode" the negative statements in MSG hashes, but in face possibly a feature in certain use cases, is that the negated statement is not clearly "readable", e.g. by direct inspection of the RDF file.

This can be a good thing, when one cares that the denied statement is not to be known by third-parties necessarily. On the other hand if the denied statements want to be made legible, one could think of add other auxiliary statements for

⁷ <http://www.dbin.org/RDFContextTools.php>

this purposes (such as the above mentioned reified N3 statements, or using agreed complementary predicate URIs modified e.g. to add "not:" in front.)

5 A Simple Semantic Web Pipe Execution Engine: Description and Implementation

Having explained the idea to encapsulate negative statements in MSG hashes and its possible benefits, we have implemented a first prototypical Semantic Web Pipe engine at the heart of the DBin 2.0 semantic web client and authoring tool implementation, which we conceive to be the basis of an effective Semantic Web application middle-ware. While DBin 0.x [9] based on a P2P infrastructure where information "flows" across peers, DBin 2.0 simply provides the user with a more controlled way to define the order and the location of the sources to import and then executes the pipes to generate a final RDF base which is then browsed and queried.

For our simple prototype, we exploit this order in evaluating RDF statements to be overridden: In the DBin piping engine, RDF sources can be either local or remote. These are ordered in a stack according to the priority selected by the user. At execution stage, a new empty triplestore is created which will contain the graph resulting from the pipe, let us call it ' T '. The sources are then processed one by one, from the one with lower priority to the one with higher priority. Naming the currently processed graph ' G ', the specialized merge procedure is the following:

1. G is cleaned by any negative MSG that overwrites a positive MSG in G (this means that if G expresses " X " and "not X " we delete both the assertions);
2. The content of G is added to T ;
3. Negative statements are "applied", so that is if positive statements exists in T corresponding to negative statements, the lower priority positive statements are removed (this step is the same of the first one except that it is applied to the resulting graph T);
4. Any remaining negative statement is dropped, as they must not have effect on the higher graphs considered in next cycles.

Once this conflict-resolving "merge" procedure has been performed for all the RDF sources, T contains the final RDF model and DBin applies RDFS reasoning on it. We remark that the result in absence of negated statements tantamounts to exactly the the common RDF merge.

Clearly, by handling conflict resolution at the RDF merge level, and applying RDFS reasoning only at the last step many issues are solved in a simple, intuitive and, at the same time, efficient manner. By removing at each step any remaining negative statement we opt for a "non symmetric" approach where positive statements are somehow considered more important and persistent than "negative" ones. Moreover, the remaining RDF set is clearly consistent (being simple RDF). We note however, that there could also be possibly problematic corner cases:

For instance, imagine that Bob sneaks in the unwanted statement about Alice as follows:

```
<http://ex.org/~charles\#me>  
  myfoaf:likes <http://alice.example.org/i>.  
myfoaf:likes rdfs:subPropertyOf foaf:knows.
```

In this disguise, even if Bob’s FOAF data is given lower priority than Charles’ FOAF file, the unwanted statement would survive the conflict resolution during our ordered merge, since we do not do RDFS inference in this process.

We are currently, investigating repairs to our approach which remedy this situation, e.g. by labeling inferred triples with the priority of the lowest statement contributing to their inference and, in a recursive process removing conflicting inferred triples in a post processing step. Unfortunately, we conjecture that finding this lowest statement is, in the general intractable⁸, but we hope that an approximative solution, which at least guarantees that only overall sound triples are inferred might be achievable.

Another drawback of the current approach is that the priority order among considered RDF sources has to be given upfront as user input to DBin, which might not be a problem for smaller scale pipe examples, but be undesirable as the number of known sources grow to large scale. Trust negotiation policies, see e.g. [3], encoded directly as RDFstatements within the sources could help to assess priorities among RDF sources as we require them directly from RDF data in those resources.

6 Conclusions and future works

We outlined in the present work a practical solution to add negative statements to RDF without generating overall logical inconsistency. Even leaving aside full OWL inference, we believe that being able to override RDF statements based on user priorities on which Web resources are trustworthy, is a crucial feature in Semantic Web applications. In this paper we first analyzed how negative statements can at all be expressed in current Semantic Web languages and came to the conclusion these languages do not properly address this problem, not providing means to override statements in a user defined priority order among RDF sources on the Web. Based on this observation, we presented a practical solution to the problem which is implemented on top of the DBin 2.0 system.

Our general ideas are based on the assumption that we believe only partially in Web scale DL reasoning handling complete OWL inferencing to be feasible in the near future. Our approach is a more practical one dealing with the increasing number of RDF data out there in an effective and arguably feasible manner. Negative statements treated in this work, which is still in a preliminary stage, are a first example of practical necessities we plan to address when effectively and efficiently processing Semantic Web data for useful Semantic Web applications in the spirit of “Semantic Web Pipes”. In this sense, this work is conceived to spark discussions for more practical solutions towards making the Semantic Web real,

⁸ A concrete algorithm and complexity studies for such an algorithm are still on our agenda

which might also raise controversy among “purists” in terms of what the term “Semantic Web Reasoning” comprises and what not. More examples of issues we want handle in practical implementations include linking RDF data by adding views⁹, possibly involving scoped negation [11, 10] and evaluate scalability of such extensions in practical scenarios.

References

1. S. Auer, C. Bizer, G. Kobilarov, J. Lehmann, R. Cyganiak, and Z. Ives. Dbpedia: A nucleus for a web of open data. In *6th Int'l Semantic Web Conference*, Busan, Korea, Nov. 2007.
2. T. Berners-Lee. Notation 3, since 1998. Available at <http://www.w3.org/DesignIssues/Notation3.html>.
3. P. A. Bonatti and D. Olmedilla. Rule-based policy representation and reasoning for the semantic web. In *Reasoning Web - Third International Summer School*, pages 240–268, Dresden, Germany, Sept. 2007.
4. J. J. Carroll. Signing rdf graphs. In *The Semantic Web - ISWC 2003, Second International Semantic Web Conference*, pages 369–384, Sanibel Island, FL, USA, Oct. 2003.
5. FOAF Vocabulary Specification, July 2005. <http://xmlns.com/foaf/0.1/>.
6. A. Harth, J. Umbrich, and S. Decker. Multicrawler: A pipelined architecture for crawling and indexing semantic web data. In *5th International Semantic Web Conference*, Athens, GA, USA, Nov. 2006.
7. P. Hayes. RDF semantics. Technical report, W3C, February 2004. W3C Recommendation.
8. Z. Huang, F. van Harmelen, and A. ten Teije. Reasoning with inconsistent ontologies. In *Proceedings of the Nineteenth International Joint Conference on Artificial Intelligence (IJCAI'05)*, Edinburgh, Scotland, August 2005.
9. M. Nucci, C. Morbidoni, and G. Tummarello. Enabling semantic web communities with dbin: an overview. In *ISWC2006 Semantic Web challenge*, Athens, GA, USA, 2006. Finalist.
10. A. Polleres, C. Feier, and A. Harth. Rules with contextually scoped negation. In *3rd European Semantic Web Conference (ESWC2006)*, volume 4011 of *Lecture Notes in Computer Science*, Budva, Montenegro, June 2006. Springer.
11. A. Polleres, F. Scharffe, and R. Schindlauer. SPARQL++ for mapping between RDF vocabularies. In *6th International Conference on Ontologies, DataBases, and Applications of Semantics (ODBASE 2007)*, Vilamoura, Algarve, Portugal, Nov. 2007. To appear.
12. E. Prud'hommeaux and A. Seaborne (eds.). SPARQL Query Language for RDF, June 2007. W3C Candidate Recommendation, available at <http://www.w3.org/TR/2007/CR-rdf-sparql-query-20070614/>.
13. G. Tummarello, R. Delbru, and E. Oren. Sindice.com: Weaving the open linked data. In *Proceedings of the International Semantic Web Conference (ISWC)*, Nov. 2007. To appear.
14. G. Tummarello, C. Morbidoni, P. Puliti, and F. Piazza. Signing individual fragments of an RDF graph. In *Special interest tracks and posters of the 14th international conference on World Wide Web*, Chiba, Japan, 2005.

⁹ see also W3C RIF Use Case 10, http://www.w3.org/TR/rif-ucr/#Publishing_Rules_for_Interlinked_Metadata