

SWI-Prolog RDF parser

Jan Wielemaker
SWI,
University of Amsterdam
The Netherlands
E-mail: `jan@swi.psy.uva.nl`

April 27, 2000

Abstract

RDF (**R**esource **D**escription **F**ormat) is a W3C standard for expressing meta-data about web-resources. It has three representations providing the same semantics. RDF documents are normally transferred as XML documents using the RDF-XML syntax. This format is very unsuitable for processing. The parser defined here converts an RDF-XML document into the triple notation.

Contents

| | | |
|----------|-----------------------------------|----------|
| 1 | Introduction | 2 |
| 2 | Parsing RDF in Prolog | 2 |
| 3 | Predicates | 3 |
| 4 | Name spaces | 4 |
| 4.1 | Low-level access | 4 |
| 5 | Testing the RDF translator | 5 |
| A | Metrics | 5 |
| B | Installation | 6 |
| B.1 | Unix systems | 6 |
| B.2 | Windows | 6 |

1 Introduction

RDF is a promising standard for representing meta-data about documents on the web as well as exchanging ontologies. RDF is often associated with ‘semantics on the web’. It consists of a formal data-model defined in terms of *triples*. In addition, a *graph* model is defined for visualisation and an XML application is defined for exchange.

‘Semantics on the web’ is often associated with the Prolog programming language. It is assumed that Prolog is a suitable vehicle to reason with the data expressed in RDF models. Most of the related web-infra structure (e.g. XML parsers, DOM implementations) are defined in Java, Perl, C or C++.

Various routes are available to the Prolog user. Low-level XML parsing is due to its nature best done in C or C++. These languages produce fast code. As XML/SGML are the basis of most of the other web-related formats we will benefit most here. XML and SGML being very stable specifications make this even more attractive.

But what about RDF? RDF-XML is defined in XML, and provided with a Prolog term representing the XML document processing it according to the RDF syntax is quick and easy in Prolog. The alternative, getting yet another library and language attached to the system, is getting less attractive.

2 Parsing RDF in Prolog

To demonstrate this, we realised an RDF compiler in Prolog on top of the sgml2pl package (providing a name-space sensitive XML parser). The transformation is realised in two passes.

The first pass rewrites the XML term into a Prolog term conveying the same information in a more friendly manner. This transformation is defined in a high-level pattern matching language defined on top of Prolog with properties similar to DCG (Definite Clause Grammar).

The source of this translation is very close to the BNF notation used by the specification, so correctness is ‘obvious’. Below is a part of the definition of RDF containers. Note that XML elements are represented using a term of the format:

```
element(Name, [AttrName = Value...], [Content ...])

memberElt(LI) ::=
    \referencedItem(LI).
memberElt(LI) ::=
    \inlineItem(LI).

referencedItem(LI) ::=
    element(\rdf(li),
           [ \resourceAttr(LI) ],
           []).

inlineItem(literal(LI)) ::=
    element(\rdf(li),
           [ \parseLiteral ],
           LI).
inlineItem(description(description, _, _, Properties)) ::=
```

```

        element(\rdf(li),
                [ \parseResource ],
                \propertyElts(Properties)).
inlineItem(LI) ::=
    element(\rdf(li),
            [],
            [\rdf_object(LI)]), !. % inlined object
inlineItem(literal(LI)) ::=
    element(\rdf(li),
            [],
            [LI]).                % string value

```

Expression in the rule that are prefixed by the `\` operator acts as invocation of another rule-set. The body-term is converted into a term where all rule-references are replaced by variables. The resulting term is matched and translation of the arguments is achieved by calling the appropriate rule. Below is the Prolog code for the **referencedItem** rule:

```

referencedItem(A, element(B, [C], [])) :-
    rdf(li, B),
    resourceAttr(A, C).

```

Additional code can be added using a notation close to the Prolog DCG notation. Here is the rule for a description, producing properties both using **propAttrs** and **propertyElts**.

```

description(description, About, BagID, Properties) ::=
    element(\rdf('Description'),
            \attrs([ \?idAboutAttr(About),
                    \?bagIdAttr(BagID)
                    | \propAttrs(PropAttrs)
                    ]),
            \propertyElts(PropElts)),
    { !, append(PropAttrs, PropElts, Properties)
    }.

```

3 Predicates

The parser is designed to operate on various environments and therefore provides interfaces at various levels. First we describe the top level defined in `library(rdf)`, simply parsing a PDF-XML file into a list of triples. Please note these are *not* asserted into the database because it is not necessarily the final format the user wishes to reason with and it is not clear how the user wants to deal with multiple RDF documents. Some options are using global URI's in one pool, in Prolog modules or using an additional argument.

```

load_rdf(+File, -Triples)
    Same as load_rdf(File, Triples, []).

```

load_rdf(*+File*, *-Triples*, *+Options*)

Read the RDF-XML file *File* and return a list of *Triples*. *Options* defines additional processing options. Currently defined options are:

base_uri(*BaseURI*)

If provided local identifiers and identifier-references are globalised using this URI. If omitted or the atom [], local identifiers are not tagged.

The *Triples* list is a list of **rdf**(*Subject*, *Predicate*, *Object*) triples. *Subject* is either a plain resource (an atom), or one of the terms **each**(*URI*) or **prefix**(*URI*) with the obvious meaning. *Predicate* is either a plain atom for explicitly non-qualified names or a term *Namespace:Name*. If *Namespace* is the defined RDF name space it is returned as the atom **rdf**. Finally, *Object* is a URI, a *Predicate* or a term of the format **literal**(*Value*) for literal values. *Value* is either a plain atom or a parsed XML term (list of atoms and elements).

4 Name spaces

XML name spaces are identified using a URI. Unfortunately various URI's are in common use to refer to RDF. The **rdf_parser.pl** module therefore defined the namespace as a **multifile/1** predicate, that can be extended by the user. For example, to parse the Netscape OpenDirectory **structure.rdf** file, the following declarations are used:

```
:- multifile
    rdf_parser:rdf_name_space/1.

rdf_parser:rdf_name_space('http://www.w3.org/TR/RDF/').
rdf_parser:rdf_name_space('http://directory.mozilla.org/rdf').
rdf_parser:rdf_name_space('http://dmoz.org/rdf').
```

The initial definition of this predicate is given below.

```
rdf_name_space('http://www.w3.org/1999/02/22-rdf-syntax-ns#').
rdf_name_space('http://www.w3.org/TR/REC-rdf-syntax').
```

4.1 Low-level access

The above defined **load_rdf**/[2,3] is not always suitable. For example, it cannot deal with documents where the RDF statement is embedded an XML document. It also cannot deal with really big documents (e.g. the Netscape OpenDirectory project), without huge amounts of memory.

For really big documents, the **sgml2pl** parser can be programmed to handle the content of a specific element (i.e. **<rdf:RDF>**) element-by-element. The parsing primitives defined in this section can be used to process these one-by-one.

xml_to_rdf(*+XML*, *+BaseURI*, *-Triples*)

Process an XML term produced by `load_structure/3` using the `dialect(xmlns)` output option. *XML* is either a complete `<rdf:RDF>` element, a list of RDF-objects (container or description) or a single description of container.

process_rdf(*+File*, *+BaseURI*, *:OnTriples*)

Exploits the call-back interface of `sgml2pl`, calling *OnTriples* with the list of triples resulting from a single top level RDF object for each RDF element in the file. This predicate can be used to process arbitrary large RDF files as the file is processed object-by-object. The example below simply asserts all triples into the database:

```
assert_list([]).
assert_list([H|T]) :-
    assert(H),
    assert_list(T).

?- process_rdf('structure,rdf', [], assert_list).
```

5 Testing the RDF translator

A test-suite and driver program are provided by `rdf_test.pl` in the source directory. To run these tests, load this file into Prolog in the distribution directory. The test files are in the directory `suite` and the proper output in `suite/ok`. Predicates provided by `rdf_test.pl`:

suite(*+N*)

Run test *N* using the file `suite/tN.rdf` and display the RDF source, the intermediate Prolog representation and the resulting triples.

passed(*+N*)

Process `suite/tN.rdf` and store the resulting triples in `suite/ok/tN.pl` for later validation by `test/0`.

test

Run all tests and classify the result.

A Metrics

It took three days to write and one to document the Prolog RDF parser. A significant part of the time was spent understanding the RDF specification.

The size of the source (including comments) is given in the table below.

| lines | words | bytes | file | function |
|-------|-------|-------|----------------------------|-------------------|
| 109 | 255 | 2663 | <code>rdf.pl</code> | Driver program |
| 312 | 649 | 6416 | <code>rdf_parser.pl</code> | 1-st phase parser |
| 246 | 752 | 5852 | <code>rdf_triple.pl</code> | 2-nd phase parser |
| 126 | 339 | 2596 | <code>rewrite.pl</code> | rule-compiler |
| 793 | 1995 | 17527 | total | |

We also compared the performance using an RDF-Schema file generated by Protege-2000 and interpreted as RDF. This file contains 162 descriptions in 50 Kbytes, resulting in 599 triples. Environment: Intel Pentium-II/450 with 384 Mbytes memory running SuSE Linux 6.3.

The parser described here requires 0.15 seconds excluding 0.13 seconds Prolog startup time to process this file. The Pro Solutions parser (written in Perl) requires 1.5 seconds excluding 0.25 seconds startup time.

B Installation

B.1 Unix systems

Installation on Unix system uses the commonly found *configure*, *make* and *make install* sequence. SWI-Prolog should be installed before building this package. If SWI-Prolog is not installed as `pl`, the environment variable `PL` must be set to the name of the SWI-Prolog executable. Installation is now accomplished using:

```
% ./configure
% make
% make install
```

This installs the Prolog library files in `$PLBASE/library`, where `$PLBASE` refers to the SWI-Prolog 'home-directory'.

B.2 Windows

Copy the files `rdf.pl`, `rdf_parse.pl`, `rdf_triple.pl` and `rewrite.pl` to the SWI-Prolog library, start Prolog and run `make/0` to update the library index.