TDT4510 Data and Information Management, Specialization Project
Fall 2008

# Developing a SPARQL parser for .NET

**Authors**:
Ole Petter Bang and Tormod Fjeldskår


**Supervisors**:
Professor Svein Erik Bratsberg and
Fast representative Øystein Torbjørnsen

NTNU



Norwegian University of Science and Technology
Department of Computer and Information Science

# Abstract

The Semantic Web is growing, both in size and popularity. At the core of the Semantic Web is the *Resource Description Framework* (RDF). RDF allows for encoding of information in web-based applications. The *SPARQL Protocol and RDF Query Language* is a SQL-like query language, recommended by W3C for querying RDF data.

A SPARQL parser front-end, called SharQL, has been created for the Microsoft .NET Framework using the *MPLex* and *MPPG* tools from Microsoft's *Managed Babel* package. The parser is capable of transforming textual SPARQL queries into abstract syntax trees (ASTs) representing the queries.

The parser has been tested using syntax tests for the SPARQL language released by the *Data Access Working Group* of W3C. The test suite contains categorized tests focusing on different aspects of the SPARQL grammar. All relevant tests pass using the SharQL parser, indicating the level of conformity with the SPARQL specification.

In order to utilize the results of the project, a back-end is needed in order to perform analysis, transformations and actual data queries against an RDF data store.

# Preface

This report is the result of the course "TDT4510 Data and Information Management, Specialization Project" at the Norwegian University of Science and Technology (NTNU), during fall 2008. It is a response to the following task description:

> RDF is a schema for encoding information in web-based applications and SPARQL is a W3C standard query language similar to SQL that is used to query such RDF data. The task of the project is to build a parser for SPARQL that will generate an abstract syntax tree (AST). The project can use open-source components but the resulting code should be written in C# and not be limited by license terms blocking commercial use.

We would like to thank our supervisors, professor Svein Erik Bratsberg at NTNU and Øystein Torbjørnsen at FAST, for feedback and guidance during the course of this project.

Trondheim, December 2008

Ole Petter Bang        Tormod Fjeldskår

# Contents

# List of Tables

# List of Figures

# Chapter 1

# Introduction

"The Semantic Web is a web of data." - W3C [3]

The *Semantic Web* attempts at a transition from the traditional web of application-centric data, to a web of interconnected data from various data sources. A corner stone of the Semantic Web is the *Resource Description Framework* (RDF). RDF allows for representation of metadata about web resources, or entities that can be identified by web resources. [4]

Describing web resources is one thing. To leverage the power of interconnected data, however, one needs a way of querying them. If e.g. a photo sharing application and a calendar sharing application both publish RDF data, it might be of interest to find all photos taken while the photographer was attending a certain photography conference.

*SPARQL* (SPARQL Protocol and RDF Query Language) is the language proposed by the *RDF Data Access Working Group* for querying RDF datasets. On January 15th 2008, the SPARQL specification reached the status of *W3C Recommendation*, the highest maturity level possible for W3C specifications. With SPARQL, queries like the aforementioned photography case are easy to express. [1]

*Information Access Disruptions* (iAD) is a constellation between FAST, two Norwegian enterprises and different research environments. Their goal is to develop the best search technology in the world. As part of this research, it is desirable to be able to query several different types of information from various data sources, including RDF.

The underlying search technology used by FAST is called *MARS*. The goal of this project is to create a parser for the SPARQL language, written in C# and running on Microsoft's .NET Framework. The SPARQL parser produced as part of this project should in turn be able to facilitate querying against RDF data sources using the .NET based MARS technology.

In Chapter 2, a brief introduction to RDF and SPARQL is given. Following this introduction, Chapter 3 presents the architectural decisions made while developing the SPARQL parser. This mainly concerns the choice of

parser generator, but also explains briefly the class hierarchy facilitating AST creation. Chapter 4 describes the steps of creating the SPARQL parser, all the way from the setup of the development environment to the encapsulation of the resulting parser. Chapter 5 explains the methods used to test the resulting parser, to ensure that it behaves according to the specification. In Chapter 6, the results of the project are presented, both in terms of how the parser operates and how well it conforms with the specification. Finally, Chapter 7 concludes and summarizes the findings of the project and presents further work necessary for this parser to facilitate iAD research.

# Chapter 2

# RDF and SPARQL

RDF is a schema for encoding information in web-based applications and SPARQL is a W3C standard query language similar to SQL that is used to query RDF data. This chapter gives a short introduction to the nature of RDF data, followed by an introduction to SPARQL.

## 2.1  RDF Essentials

RDF (Resource Description Framework [4]) data represents information (particularly metadata) about resources in the World Wide Web. The concept of a "resource" is typically generalized, allowing for information not directly retrievable on the Web to be represented. Thus, it is entirely up to the producers and consumers of the RDF data to agree upon the semantics of the information, as long as it is represented according to the RDF schema. Different efforts have been made in order to standardize such semantics, including The Dublin Core Initiative [5].

RDF is intended for situations in which information needs to be processed by applications and exchanged without loss of meaning. Resources are identified using URIs[1], and described in terms of properties and their corresponding values. The properties are also identified using URIs, and their values are identified either by URIs or data typed literals. Altogether, resources and their properties are represented as *triples*, consisting of a subject (the resource), a predicate (the property), and an object (the property value). This enables RDF to represent simple statements about resources as a *graph* of nodes with arcs representing the resources, their properties and corresponding values.

---

[1]Actually, the term URI in the context of both RDF and SPARQL always refers to Internationalized Resource Identifiers (IRIs), a generalization of the Uniform Resource Identifier (URI), which may contain characters form the Universal Character Set (Unicode/ISO 10646).

Figure 2.1: RDF sample (graph)

```
<rdf:RDF
    xmlns:rdf="http://www.w3.org/1999/02/22-rdf-syntax-ns#"
    xmlns:dc="http://purl.org/dc/elements/1.1/">
  <rdf:Description
      rdf:about="http://en.wikipedia.org/wiki/Tony_Benn">
    <dc:title>Tony Benn</dc:title>
    <dc:publisher>Wikipedia</dc:publisher>
  </rdf:Description>
</rdf:RDF>
```

Figure 2.2: RDF sample (XML)

Figure 2.1 illustrates a sample RDF graph. The resource in question is a Wikipedia web-page about the British politician Tony Benn, represented by the elliptic node labeled with the web-page URI. The standardized properties used to describe the resource are gathered from Dublin Core [5], each represented by an arc pointing at the corresponding rectangular property value node.

The XML-based representation of the very same information is given in Figure 2.2. A few, just as common, non-XML-based representations exist as well, like the Turtle [6] representation illustrated in Figure 2.3. The most obvious advantage non-XML-based representations like Turtle have over XML-based representations, is the strongly reduced level of verbosity. Also, both the Turtle and SPARQL languages are subsets of the Notation 3 [7] language, making the two syntactically very much alike.

```
@prefix dc: <http://purl.org/dc/elements/1.1/>.

<http://en.wikipedia.org/wiki/Tony_Benn>
  dc:title "Tony Benn";
  dc:publisher "Wikipedia".
```

Figure 2.3: RDF sample (Turtle)

## 2.2 Introduction to SPARQL

SPARQL (SPARQL Protocol and RDF Query Language [1]) is a query language designed to meet the use cases and requirements identified by the RDF Data Access Working Group [8].

SPARQL queries may consist of triple patterns, conjunctions, disjunctions and optional patterns. Most forms of queries contain a set of triple patterns, which are just like RDF triples except that each of the subject, predicate and object may be a variable. Such triple patterns match an RDF subgraph when an equivalent graph may be constructed from the triple pattern by substituting all variables with the corresponding terms from the subgraph.

Given the RDF data illustrated in Figure 2.2, the SPARQL query illustrated in Figure 2.4 matches the RDF subgraph describing the title property. This is done by substituting the *?title* variable with the property value "Tony Benn", which is also bound to the variable and then returned as illustrated in Table 2.1.

```
PREFIX dc: <http://purl.org/dc/elements/1.1/>

SELECT ?title
WHERE
{
  <http://en.wikipedia.org/wiki/Tony_Benn> dc:title ?title
}
```

Figure 2.4: SPARQL sample query

| title |
|-------|
| Tony Benn |

Table 2.1: SPARQL sample query result set

Each solution gives one way in which the selected variables can be bound to an RDF subgraph so that the query pattern matches the data. When literals are matched, the data type is also considered. In addition to common data

types, strings may be suffixed with a language tag, and custom data types may be used.

To further restrict the matches, *filters* can be used to restrict common data type values as well as to test different kinds of conditions using test functions. Restrictions on common data types include regular expressions on strings and arithmetic expressions for numeric values.

In an RDF triple the subject and object may be *blank nodes*, which are nodes without any identification other than a label to differentiate nodes from each other. Generally, blank nodes in a result set are not required to have the same labels as their corresponding blank nodes from the SPARQL query. The blank nodes simply act as non-distinguished variables in the query, not as references to specific blank nodes. Thus, an application writer should not expect a blank node in a query to refer to a particular blank node in the data.

Specific triple patterns may be marked as being optional, allowing for solutions to exist where the optional triple patterns does not match any RDF subgraph. Further, alternative triple patterns may be given, allowing for only one of several alternative graph patterns to match. If more than one of the alternatives match, all the possible pattern solutions are found. Also, several RDF graphs may be included and merged in a query, to allow query of several RDF stores.

The solution sequence of a SPARQL query may be influenced by a set of modifiers, most of which are familiar from languages like SQL. The modifiers available allow for ordering, projecting, selecting distinct solutions only, setting an offset defining the element position in the solution sequence from which to retrieve elements, as well as limiting the number of elements to be retrieved, starting from the beginning of the solution sequence unless an offset is stated explicitly. In addition a *REDUCED* modifier exists. The cardinality of any set of variable bindings in a reduced solution set is at least one and not more than the cardinality of the solution set with no *DISTINCT* or *REDUCED* modifiers.

SPARQL has four query forms. *SELECT* queries return all, or a subset of, the variables bound in a query pattern match. *CONSTRUCT* queries return RDF graphs constructed by substituting variables in a set of triple templates. An *ASK* query returns a boolean value indicating whether a specific query pattern matches or not. *DESCRIBE* queries return RDF graphs describing the resources found.

## 2.3   SPARQL Testimonials

Many companies have taken an interest in SPARQL. Following the press release announcing the publication of SPARQL [9], several early adopters responded with testimonials about W3C's new specification [10]. Among

these were both smaller and larger companies. To give an impression of how SPARQL is received in the communities, some of the testimonials are quoted below.

**Computas AS:**

>    *Computas AS is currently building systems for its customers where SPARQL is a fundamental core component. When conducting feasibility studies, we found that there are allready many high quality off-the-shelf components that puts the vision of the data web within reach, also for smaller enterprises.*
>
>    *We are pleased to see SPARQL promoted to a W3C Recommendation, as it provides a stable platform for further work. We are allready experimenting with extensions to SPARQL, and will work with the W3C and its membership in the work that lies ahead.*
>
>    **– Kjetil Kjernsmo, Senior Knowledge Engineer, Computas AS**

**Hewlett-Packard:**

>    *Hewlett-Packard is pleased to support the SPARQL Recommendations.*
>
>    *SPARQL is a key element for integrated information access across information silos and across business boundaries. HP customers can benefit from better information utilization by employing semantic web technologies.*
>
>    *HP's Jena Semantic Web framework has a complete implementation of query language, protocol and result set processing. Jena is open-source, freely available, with a large and active developer community.*
>
>    *HP is pleased to announce the first full release of SDB, a new SPARQL database system for Jena that leverages existing database installations to give enterprise-grade storage and query of RDF.*
>
>    **– Jean-Luc Chatelain, CTO HP Software Information Management**

**Oracle:**

>    *Oracle congratulates the W3C on achieving 'Recommendation' status for SPARQL. As an active participant in this working group, Oracle believes the standardization of SPARQL will play an instrumental role in achieving the vision of the Semantic Web. The community's work is intended to help organizations more effectively discover, automate, integrate and re-use data across various applications.*

*Oracle Database 11g Semantic Store provides native support for efficient and scalable storage, bulk loading, inferencing, and graph-pattern based querying of semantic data represented using W3C's RDF, RDFS, and OWL languages. The Oracle Jena adaptor allows querying of semantic data stored in Oracle using the SPARQL query language while leveraging the performance and scalability of Oracle's Semantic Store.*

**– Don Deutsch, vice president Standards Strategy and Architecture, Oracle**

# Chapter 3

# Architectural Decisions

This chapter presents the architectural decisions made. This mainly involves choosing a suitable parser generator, but the class hierarchy facilitating AST creation is also explained. Various alternatives are presented along with the rationale behind the final decisions.

## 3.1 Choosing a Parser Generator

The main architectural decision for this project was the choice of parser generator. A parser generator is a program that reads a language specification formatted using a certain syntax. Based on this specification, a parser is created. Different generators create different parser types. Some create top-down (recursive descent) *LL(k)* parsers, while others create bottom-up *LALR* parsers. LALR parsers can handle a wider range of grammars than LL(k) parsers [11].

Given that the resulting parser of this project is to be written in C#, the parser generator has to be able to generate C# source code. According to the SPARQL specification, *"The SPARQL grammar is LL(1) when the rules with uppercased names are used as terminals."* [1]. This section discusses some of the alternatives that were considered, and presents the rationale behind the choice.

### 3.1.1 ANTLR

*ANTLR* (<u>AN</u>other <u>T</u>ool for <u>L</u>anguage <u>R</u>ecognition) is an open source parser generator, written in Java. It is considered mature and has reached version 3.1 at the time of writing.

The parsers created by ANTLR are recursive descent LL(k) parsers which are less expressive than LR/LALR parsers, but the generated source code is more intuitive. Source code can be generated in several programming

|  | **ANTLR** | **Coco/R** | **MPLex/MPPG** |
|---|---|---|---|
| Rules definition | EBNF | Attributed EBNF | Yacc-like BNF |
| Parser type | LL(k) | LL(k) | LALR |
| AST generation | Automatic | Manual | Manual |
| Supported by | Community | Community | Microsoft |
| Visual Studio integration out-of-the-box | No | No | Yes (requires manual editing of project file) |

Table 3.1: Summary of parser generators

languages, such as Java, C++ and C#. Rules are defined in a format similar to EBNF. [12]

### 3.1.2 Coco/R

*Coco/R* (<u>Co</u>mpiler <u>c</u>ompiler/<u>R</u>ecursive descent) is similar to ANTLR in that is creates recursive descent LL(k) parsers. However, while ANTLR is written in Java with the possibility of generating source code in various languages, Coco/R comes in different versions supporting different languages. To generate C# code, the C# version of Coco/R is required. Rules are defined as an attributed EBNF grammar [13].

### 3.1.3 MPLex and MPPG

*MPLex* (<u>M</u>anaged <u>P</u>ackage <u>Lex</u>) and *MPPG* (<u>M</u>anaged <u>P</u>ackage <u>P</u>arser <u>G</u>enerator) are provided by Microsoft as part of their Visual Studio SDK. The scanner generator *MPLex* and the parser generator *MPPG* are part of the Managed Babel package aimed at Visual Studio extension developers. [14]

MPLex and MPPG are closely related to the *Queensland University of Technology* open source projects *Gardens Point Scanner Generator* (GPLEX) [15] and *Gardens Point Parser Generator* (GPPG) [16]. The parsers generated by MPPG are C# only, bottom-up LALR parsers.

Being part of the Visual Studio SDK, integrating the actions of MPLex and MPPG into a Visual Studio project is relatively easy. This is an advantage over GPLEX and GPPG. Rules for MPLex and MPPG are specified in a format largely based on the syntax used by tools like Lex and Yacc [17].

10

### 3.1.4 Our Choice

The various properties of the different parser generators are listed in Table 3.1. In cooperation with the supervisors, the choice of parser generator ultimately fell on MPLex/MPPG. The LALR parsers created by MPPG are more expressive than the LL(k) parsers created by the two other alternatives [11]. This could be a future advantage if changes to the SPARQL language require such expressiveness. Even though ANTLR and Coco/R have large communities and good testimonials, the fact that MPLex and MPPG are provided and supported by Microsoft [14] was decisive. The ease of integration with the Visual Studio environment is an advantage as well.

## 3.2 Further Architectural Choices

As the actual code performing the lexical analysis and the parsing itself is generated using MPLex and MPPG, the remaining functionality subject to architectural decisions is the construction and representation of the abstract syntax tree (AST) representing the parsed queries.

If the parser was to be generated using Lex and Yacc [17], before the entry of object oriented languages, the abstract syntax tree would most likely have been constructed and represented using C variable pointers to *structs*, imitating objects. The C programming language has no knowledge of objects, making this approach the only viable option for implementations of notable size and complexity.

Inspired by its predecessors, the C# programming language still retains knowledge of structs. C# variables always represent structs by value; a struct may never be represented by reference, though it may be passed by reference between methods [18]. When reassigning a variable holding a C# struct, the entire content of the struct is copied to the memory location references by the variable. A variable referencing a class object would instead have changed the destination of the reference to where the class object already resides. Also, C# structs can only implement interfaces, they cannot inherit other classes or structs. Not being able to exploit the powerful concept of inheritance makes C# structs even less suitable for AST representation.

An AST is by nature based on references between tree nodes. Thus class objects are the only viable option for representing AST nodes. Using class objects, a suitable class hierarchy may be built, benefiting from important properties of object orientation like inheritance and shared interfaces in order to generalize the hierarchy.

For this project, the planned architecture involves using C# objects instantiated to represent the abstract syntax tree. Common base functionality should be implemented in a single base node, from which all other nodes

inherit. This is further discussed in Section 4.2.

# Chapter 4

# Implementation

This chapter describes the steps of creating the SPARQL parser. First, the setup of the development environment is explained. Then the process of translating W3C's SPARQL specification to a format that MPLex and MPPG understand is described. Error handling, along with some necessary pre- and post-processing is then explained, as well as how the entire implementation is encapsulated in a facade class. Finally, a test client for visual inspection of resulting ASTs is presented.

## 4.1 Setup of the Development Environment

MPLex and MPPG are designed to easily integrate with Microsoft Visual Studio 2008. A Visual Studio solution can be configured to automatically generate source files for the scanner and parser, and then compile them along with the rest of the source files.

Before configuring the development environment, the following software had to be installed.

**Microsoft Visual Studio 2008**  The IDE used during the entire development, provided via MSDN Academic Alliance [19]. A 90-day trial version is available from `http://msdn.microsoft.com/en-us/vstudio/products/aa700831.aspx`.

**Visual Studio 2008 SDK 1.1**  Contains the tools MPLex and MPPG. Available from `http://www.microsoft.com/downloads/details.aspx?FamilyID=59ec6ec3-4273-48a3-ba25-dc925a45584d`.

The setup of a Visual Studio solution for autogeneration of the parser involves manual editing of the corresponding *project file*. A project file in Visual Studio is similar in semantics to a makefile in Unix development. It describes the various sub-actions performed as part of the main build action.

Figure 4.1: Creating the Visual Studio project

First, the Visual Studio solution with the parser project was created by clicking *File → New → Project...* and then choosing a C# Class Library, as shown in Figure 4.1.

The autogenerated file *Class1.cs* was deleted, and two new text files were added, called *Scanner.lex* and *Parser.y*. These two files were the ones that would later contain the SPARQL language specifications and dictate the generation of the parser. In order for these files to be handled properly by MPLex and MPPG, their encoding had to be set by clicking *File → Advanced Save Options...*. In the dialog that appeared, *Western European (Windows) - Codepage 1252* had to be set as the encoding as shown in Figure 4.2.

The scanner also requires an interface called *IErrorHandler* in order to communicate any irregularities in the input. Thus, an empty interface *IErrorHandler* was created and assigned public visibility. The solution layout now looked like shown in Figure 4.3.

At this point, the solution would build, but no parser would be generated. In order to set up MPLex and MPPG to process *Scanner.lex* and *Parser.y*, the project file had to be edited manually. To edit the project file, the project must first be unloaded. This was done by right-clicking the project root and selecting *Unload Project*. With the project unloaded, the project file can be edited by right-clicking the project root again and selecting *Edit [ProjectName].csproj* as shown in Figure 4.4.

14

Figure 4.2: Choosing the right encoding



Figure 4.3: Visual Studio solution after adding initial items

Figure 4.4: Open the Visual Studio project file for editing

Towards the end of the project file, an ItemGroup element contained the scanner and parser specification files. Because both Parser.y and Scanner.lex was placed inside None-elements, no action would be applied to them. To let MPLex and MPPG handle these files, *Parser.y* had to be placed in an *MPPGCompile* element and *Scanner.lex* in an *MPLexCompile* element, as shown in Figure 4.5. When reloading the project by right-clicking the project root and selecting *Reload Project*, the changes made were reflected in the Properties explorer as shown in Figure 4.6; for *Parser.y*, the build action was set to *MPPGCompile* and for *Scanner.lex* it was set to *MPLexCompile*.

Also at this point, the solution would build, but no parser was generated. Setting the appropriate build actions was not sufficient for Visual Stu-



Figure 4.5: Actions for MPLex and MPPG

16

Figure 4.6: Build actions reflected in Properties explorer

```
%namespace SharQL

%%
```

Figure 4.7: An empty scanner declaration

dio to perform the actual parser generation. The assembly *Microsoft.VsSDK.Build.Tasks.dll* from the Visual Studio SDK was needed, and had to be copied into the project. A *lib* folder was created as part of the project. The assembly was added by right-clicking the new folder and selecting *Add → New Item...*. The file is usually located in the folder *%PROGRAMFILES%\MSBuild\Microsoft\VisualStudio\v9.0\VSSDK*.

The final step in setting up the integration of MPLex and MPPG with the Visual Studio solution was to hook up the project file to the newly added assembly. This was done by opening the project file for manual editing once again and pasting the XML from Appendix A immediately before the closing tag of the *Project* element.

Reloading the project and attempting to build the project at this point yielded an error message from MPLex: "*Parser error <syntax error, unexpected EOF>*". This was good news, as MPLex was obviously trying to parse the *Scanner.lex* file, but was unable to do so because the file was empty. Adding the code from Figure 4.7 to *Scanner.lex* fixed the error. At the same time, the code from Figure 4.8 was added to *Parser.y* to ensure a (quite useless) parser would be generated.

While this fixed the syntax error of *Scanner.lex*, three new errors arose, this time in the file *Scanner.cs*. These errors were all due to missing references. To resolve this issue, two more files had to be added to the *lib* folder of the project. *MPLex.exe* and *MPPG.exe* were copied from the *[SDK Root]\VisualStudioIntegration\Tools\Bin* folder. These are the tools

```
%namespace SharQL

%start Query

%%

Query :;

%%
```

Figure 4.8: An empty parser declaration

that generate the scanner and the parser, but they also contain some of the data types necessary for the scanner and parser to operate. With these files copied to the *lib* folder, a reference was added to each of them by right-clicking the *Reference* folder and selecting *Add Reference...*. On the *Browse* tab, MPLex.exe and MPPG.exe was added as shown in Figure 4.9.

With these references resolved, the solution was now building properly. Source code for both a scanner and a parser was automatically generated, resulting in a .dll-file readily compiled to the output directory. The generated scanner and parser source code is temporarily stored in the *obj* folder during a build, and is left there for later inspection if desired.

## 4.2 AST Preparations

MPLex and MPPG have no notion of an abstract syntax tree per se. In order for the resulting parser to produce an AST, node objects have to be created and structured as part of the semantic actions of the grammar productions. Node objects are created from node classes, and these classes constitute the AST class hierarchy. All the entities of this hierarchy reside in the *SharQL.Ast* namespace.

### 4.2.1 Designing the AST Class Hierarchy

At the root of the AST class hierarchy is a simple interface, *INode*, shown in Figure 4.10. This interface supports the construction of a tree structure, navigable in both directions through the *Parent* and *Children* properties. The *Accept* method is part of the *Visitor* interface, explained in Section 4.4.

Below the *INode* interface in the class hierarchy is the abstract class *NodeBase*, shown in Figure 4.11. The complete implementation of *NodeBase* is provided in Appendix D. The purpose of this class is to provide a basic implementation of the *INode* interface. All these implementations are virtual, allowing for subclasses to override them if necessary. Although this

Figure 4.9: Adding necessary references



Figure 4.10: The *INode* interface

Figure 4.11: The *NodeBase* abstract class

class has no abstract members, it is marked as abstract to prevent instantiation, and to allow for future abstract members.

The public *Parent* and *Children* properties expose the protected *parent* and *children* fields respectively. The *NodeType* read-only property is a unique identifier for each class, used as part of the Visitor pattern. The basic implementation uses reflection to retrieve the full name of the class. This property should be overridden by subclasses if performance is critical.

*NodeBase* provides two constructors. The default parameterless constructor makes sure the *children* field is initialized as an empty collection, while a second constructor allows for initializing the *children* field with initial *NodeBase* objects.

The *Children* property of *NodeBase* is exposed as an *IList<INode>* interface, and implemented as a custom *NodeCollection* class. This class is marked as *internal* and is thus not exposed by the API. A custom *IList* implementation was chosen over the standard ones to allow for customization.

The concrete AST classes all reside in the *SharQL.Ast.Nodes* namespace. Ultimately, they are all descendants of *NodeBase* and thus implement the *INode* interface.

```
%union { public NodeBase Value; }
```

Figure 4.12: This project's *%union* specification

```
Prologue : BaseDecl PrefixDeclList
{
  $$.Value = new PrologueNode($1.Value, $2.Value);
}
```

Figure 4.13: Building the AST via *ValueType*

## 4.2.2 Working With MPPG's *ValueType*

MPPG creates a struct called *ValueType* which is used for the implicit objects
($$, $1, $2 etc.) available in the grammar productions. The *%union* construct is used in the parser specification file to specify which fields should be available in the *ValueType* struct. In this project, this *%union* construct is specified as shown in Figure 4.12, adding one single field: our AST node base type.

Without any further modifications, the AST can now be constructed as part of the semantic specification for all the productions, as shown in Figure 4.13.

Because our *ValueType* contains only one field, it would make the code more readable if the implicit objects of each production represented *NodeBase* objects rather than *ValueType* objects. Using the *operator overloading* feature of C#, it is possible to allow implicit casting between the *ValueType* and *NodeBase* types. This is done as shown in Figure 4.14. Note that *ValueType* is defined as a *partial struct*, making it trivial to extend.

With implicit casting enabled between *ValueType* and *NodeBase*, it is possible to assign a *NodeBase* object to a *ValueType* variable and vice versa.

```
public partial struct ValueType
{
  public static implicit operator NodeBase(ValueType t)
  {
    return t.Value;
  }
  public static implicit operator ValueType(NodeBase t)
  {
    return new ValueType() { Value = t };
  }
}
```

Figure 4.14: Implicit casting between *ValueType* and *NodeBase*

```
Prologue : BaseDecl PrefixDeclList
{
  $$ = new PrologueNode($1, $2);
}
```

Figure 4.15: Building the AST using implicit casting

```
list1           ::= list_item*
list2           ::= list_item+
list_item       ::= OPTIONAL_PART? REQUIRED_PART
```

Figure 4.16: EBNF sample rules

Thus, the production in Figure 4.13 could be rewritten as shown in Figure 4.15.

## 4.3   Specifying the Grammar

W3C specifies the SPARQL grammar using an *Extended Backus-Naur Form* (EBNF) based notation [20]. Since the MPLex and MPPG scanner and parser generators are two separate generators, the grammar must be split into a scanner part and a parser part. Further, the generators do not support specifying grammars directly using EBNF. Instead the grammar must be specified using specification languages closely resembling the BNF-based Lex and Yacc [17] specification languages.

These specification languages lack many of the extensions of EBNF, requiring several measures to be taken in the translation of the grammar.

### 4.3.1   EBNF and BNF

EBNF is an extended version of BNF, notably introducing the modifiers **?** (zero-or-one), **\*** (zero-or-more) and **+** (one-or-more). Consider the rules in Figure 4.16. Two types of lists are specified, one that can be empty, and one that can not be empty. Further, list items have an optional part and a required part.

BNF lacks the aforementioned modifiers. Thus, to describe the same grammar using BNF, some translation is needed. Figure 4.17 shows how this is typically done ($\epsilon$ denoting the empty string). The **\*** and **+** modifier from EBNF are achieved recursively in BNF, shown in the *list_tail* rule. The **?** modifier is achieved either by enumerating all legal combinations of required and optional elements, or by introducing a new non-terminal which may produce the optional part or the empty string.

22

```
list1          ::= list_tail
list2          ::= list_item list_tail
list_tail      ::= list_item list_tail | ε
list_item      ::= REQUIRED_PART | OPTIONAL_PART REQUIRED_PART

alternatively :
list_item      ::= optional_part REQUIRED_PART
optional_part ::= OPTIONAL_PART | ε
```

Figure 4.17: BNF translation from Figure 4.16

## 4.3.2 The Scanner Specification

The grammar tokens are defined in the scanner specification using regular expressions. Each token that should be returned to the parser also defines a belonging snippet of code for returning the corresponding token enumeration item recognized by the parser. Additionally, helper tokens used to build other tokens may be defined, to avoid repeating regular expressions.

### Order of declaration

Token declaration order is of significance and must be taken into account. During token matching, the scanner will attempt to match the longest token possible. If there is a tie between two or more tokens, the one defined first is returned to the parser.

For the SPARQL grammar, however, the declaration order is not an issue. All tokens sharing a common prefix differ in length, thus a tie will never occur.

## 4.3.3 The Parser Specification

The grammar rules are declared in the parser specification. A rule defines one or more belonging productions, which are either non-empty defining a belonging snippet of code returning a token value or AST node to its superior production, or empty allowing for the production to be absent in the input string.

### Grammar Constructs

The modifiers *, + and ? are absent in the BNF-based parser specification language, and must be realized as discussed in Section 4.3.1. Figure 4.18 shows a production from the SPARQL grammar specification containing a conditional grammar symbol, namely the *WHERE* keyword. In this specific case an additional production not specifying the *WHERE* keyword is

23

```
[13] WhereClause    ::=   'WHERE'? GroupGraphPattern

/* [13] */
WhereClause
: WHERE GroupGraphPattern
{
  $$ = new WhereClauseNode($2);
}
| GroupGraphPattern
{
  $$ = new WhereClauseNode($1);
}
;
```

Figure 4.18: SPARQL Grammar Conditional Symbol Sample

specified. The bracketed numbers correspond to the rule numbers in the W3C SPARQL grammar specification.

Declaration of lists of grammar symbols (one-or-more and zero-or-more instances) requires the introduction of a replacement rule, allowing for a recursive list of grammar symbols to be constructed during parsing. Figure 4.19 shows the realization of the *DataSetClause** zero-or-more instances list using the recursive list rule *DataSetClauseList*. The rule may result in an empty production (zero instances) or a production consisting of the non-terminals *DataSetClause* and *DataSetClauseList* (more instances).

**Conflicts**

MPPG creates shift-reduce parsers, consisting of a stack holding grammar symbols and an input buffer holding the rest of the input string to be parsed. The parser shifts input symbols onto the stack until it is ready to reduce a string of grammar symbols on the top of the stack into a superior grammar production. This process is repeated until an error is detected or until the stack contains the predefined grammar start symbol and the input is empty. [11]

Two types of conflicts may occur during shift-reduce parsing: shift/reduce conflicts and reduce/reduce conflicts. Shift/reduce conflicts occur when the parser cannot decide whether to shift input symbols onto the stack or to reduce grammar symbols on top of the stack. Reduce/reduce conflicts occur when the parser cannot decide which superior production to reduce the grammar symbols on top of the stack into. The former conflict typically occurs as a result of an ambiguous grammar, and is critical. The latter easily occurs when creating the parser specification from a grammar specification like the SPARQL grammar specification, because of the required rewriting from an EBNF-based grammar to a BNF-based gram-

```
[6] ConstructQuery    ::=   'CONSTRUCT' ConstructTemplate
                            DatasetClause* WhereClause
                            SolutionModifier

/* [6] */
ConstructQuery
: CONSTRUCT ConstructTemplate DatasetClauseList WhereClause
  SolutionModifier
{
  $$ = new ConstructQueryNode($2, $3, $4, $5);
}
;


DatasetClauseList
: DatasetClause DatasetClauseList
{
  $$ = new DatasetClauseListNode($1, $2);
}
| /* empty */
{
  $$ = null;
  Debug.WriteLine("Empty DataSetClauseList");
}
;
```

Figure 4.19: SPARQL Grammar List Symbol Sample

```
/* [2] */
Prologue
: BaseDecl
{
  $$ = new PrologueNode($1, null);
}
| BaseDecl PrefixDeclList
{
  $$ = new PrologueNode($1, $2);
}
| PrefixDeclList
{
  $$ = new PrologueNode(null, $1);
}
|
{ /* empty */ }
;

PrefixDeclList
: PrefixDecl PrefixDeclList
{
  $$ = new PrefixDeclListNode($1, $2);
}
|
{ /* empty */ }
;
```

Figure 4.20: SPARQL Grammar Reduce/Reduce Sample

mar.

Figure 4.20 addresses a reduce/reduce conflict encountered in translating the SPARQL grammar specification into parser grammar productions. In this specific case, both the *Prologue* and *PrefixDeclList* rules contain an empty production. In the context of a *Prologue* grammar rule, whenever the parser encounters an empty input string, it cannot decide whether to reduce the empty production to a *PrefixDeclList* rule or to a *Prologue* rule.

This reduce/reduce conflict may be prevented simply by removing the empty production belonging to the *Prologue* rule, since a *Prologue* rule may still produce an empty production via the *PrefixDeclList* rule.

## 4.4 The Visitor Pattern

The basic idea of the *Visitor* pattern is to separate the algorithm from the data structure. The Visitor pattern allows for defining new operations on the elements of an object structure, without changing the classes of the elements on which it operates. [21]

The *Visitor* pattern is commonly used to traverse and perform opera-

26

Figure 4.21: The *IParserVisitor* interface

```
public virtual void Accept(IParserVisitor visitor)
{
  foreach (INode child in children)
  {
    if (child != null)
      child.Accept(visitor);
  }
  visitor[NodeType](this);
}
```

Figure 4.22: Basic implementation of visitor pattern from *NodeBase*

tions on abstract syntax trees. Each node in the ASTs produced by SharQL has an *Accept* method which takes a single *IParserVisitor* object. The IParserVisitor interface is shown in Figure 4.21.

The only member of this interface is an indexer which takes a string identifier as input and returns an *Action<INode>* delegate. A delegate is essentially a type-safe method pointer, and the Action<INode> delegate may point to any method which accepts a single *INode* argument and returns *void*. The AST node that accepts a visitor will access this indexer, provide its *NodeType* string as the identifier and invoke the delegate it gets in return, providing a reference to itself. An implementation snippet of this is shown in Figure 4.22. Note that the visitor will traverse the tree in a depth-first manner; before being applied to a node, it is applied to any child nodes.

### 4.4.1 Visitor Example

Consider a sample AST node as shown in Figure 4.23. A visitor implementation that traverses the AST and performs a task on all instances of this class is shown in Figure 4.24.

```
public class SampleNode : NodeBase
{
  public static readonly string SampleNodeType
    = "SampleNodeIdentifier";

  public override string NodeType
  {
    get { return SampleNode.SampleNodeType; }
  }
}
```

Figure 4.23: Sample AST node

```
public class ParserVisitor : IParserVisitor
{
  public Action<INode> this[string nodeType]
  {
    get
    {
      if (nodeType == SampleNode.SampleNodeType)
        return VisitSampleNode;
      else
        return NoOp;
    }
  }

  public void VisitSampleNode(INode node)
  {
    SampleNode n = node as SampleNode;
    if (n != null)
    {
    // Do something
    }
  }

  public void NoOp(INode node)
  {
  // Ignore unrecognized node
  }
}
```

Figure 4.24: Sample visitor implementation

## 4.5 Error Handling

The purpose of error handling in the context of the SharQL parser, is identification and reporting of error conditions concerning the syntactical correctness of a SPARQL input query. The scanner basically looks for defined tokens in the input stream and passes them on to the parser, which in turn puts a set of tokens in the context of a defined production rule. On the occurrence of an undefined, missing or unexpected token, an error condition is present.

### 4.5.1 Traditional Parser Error Handling

In a traditional scanner-parser combination like Lex and Yacc [17], either the scanner immediately reports an error on the occurrence of an unexpected input stream character, or the character is passed on to the parser, possibly in the form of a special error token, which is in turn reported as unexpected by the parser. As the scanner only knows what characters to expect on a per-token basis, the informational value of a scanner error report is limited compared to a parser error report. The parser not only knows which unexpected token that was encountered, but reports back what it expected as well.

Combining token location information from the scanner with token context information from the parser would offer an improved informational value for the error reports, and is the approach taken for the SharQL parser. Achieving such error handling requires some further extensions to both the scanner and parser.

### 4.5.2 Adapting the Grammar Definition

Section 4.3 describes the specification of the scanner and parser grammars. To support error handling, the scanner grammar definition must be extended further.

After translating the W3C SPARQL grammar, undefined tokens encountered during scanning of the input stream were simply ignored, and were not reported to the parser. Figure 4.25 shows the additional token definition needed for handling identification of undefined tokens. The token definition simply matches any single character, and returns the character itself to the parser. It is important to place the token definition below all other regular token definitions in the scanner grammar definition. This is the only way to have all valid single-character token definitions remain unaffected, as described in Section 4.3.2.

```
.               { return (int)yytext[0]; }
```

Figure 4.25: Scanner grammar definition extension

```
%YYLTYPE CustomLexLocation
```

Figure 4.26: Specifying the Token Location Type

### 4.5.3 Tracking Token Locations

When the scanner matches tokens from the input stream, token location information is available for every token matched. Even undefined tokens will get matched thanks to the adaption of the grammar as described in Section 4.5.2. This means that location information will be available for each and every token, defined or undefined, passed on to the parser.

Location information is available through global variables in the context of the scanner. The information of interest is stored in the integer variables *yyline*, *yycol* and *yyleng*. These values represent the *line*, *column* and *length*, respectively, of the currently identified token.

A special variable *yylloc* also exists, that holds all the location information of the current token. When this variable is set, it may be retrieved by the parser for further use. The type of this variable is defined in the parser specification using the *%YYLTYPE* operator, as shown in Figure 4.26. The type used by SharQL for storing the token location information is *CustomLexLocation*, a custom class replacing the *LexLocation* class supplied with MPPG. The only reason a custom class is needed is due to a bug in the *LexLocation* class occuring when two instances are merged using the belonging *Merge* method. Our *CustomLexLocation* class simply corrects this bug, and otherwise acts just as the *LexLocation* class would, as shown in Figure 4.27. When the parser reduces tokens on the stack into a production, the *CustomLexLocation* objects for all the tokens involved are merged to represent the location information for the production as a whole.

The scanner specification is altered to call a *LoadYylval* method every time a token is identified. This method constructs an instance of the *CustomLexLocation* class, supplying it with the necessary location information available from the scanner, and assigns it to the *yylloc* variable, making it available to the parser, as shown in Figure 4.28.

### 4.5.4 Collecting Errors

When the parser encounters an unexpected token or when an expected token is missing, it calls the *yyerror* method on the scanner instance. An empty *yyerror* method is already declared, which is overridden in the scan-

```
public class CustomLexLocation : IMerge<CustomLexLocation>
{
  public int sLin; // Start line
  public int sCol; // Start column
  public int eLin; // End line
  public int eCol; // End column

  public CustomLexLocation()
  {
  }

  public CustomLexLocation(int sl, int sc, int el, int ec)
  {
    sLin=sl;
    sCol=sc;
    eLin=el;
    eCol=ec;
  }

  public CustomLexLocation Merge(CustomLexLocation last)
  {
    // This part is missing from the MPPG LexLocation class
    if (last == null)
    {
      return this;
    }

    return new CustomLexLocation(sLin, sCol, last.eLin,
      last.eCol);
  }
}
```

Figure 4.27: *CustomLexLocation* class replacing the *LexLocation* class.

```
// Called each time a token has been returned.
internal void LoadYylval()
{
  // Collects token location information and makes it available
  // for the parser through the yylloc variable.
  yylloc = new CustomLexLocation(yyline, yycol, yyline, yycol
    + yyleng);
}
```

Figure 4.28: Tracking Token Location

```
// Called by the parser on the occurrence of an unexpected or
// missing token.
public override void yyerror(string s, params object[] a)
{
  // Adds the reported error to the error handler.
  handler.AddError(s, new MPLEX.Parser.LexSpan(tokLin, tokCol,
    tokLin, tokECol, tokPos, tokEPos, new ScanBuffProxy(buffer)));
}
```

Figure 4.29: *yyerror* method.

ner specification. This method is used to collect the errors reported.

The generated scanner provides a property for supplying an error handler, which must implement an *IErrorHandler* interface. The interface itself is not defined, it is just referenced. Thus, an *IErrorHandler* interface has been defined, as decribed in Section 4.1.

MPPG supplies an *ErrorHandler* class providing basic error handling. For some reason, this class does not implement the *IErrorHandler* interface referenced by the scanner. In order to re-use the existing functionality in compliance with the scanner, a custom *SharQLErrorHandler* class inheriting the *ErrorHandler* class and implementing the *IErrorHandler* interface is used.

Figure 4.29 shows the *yyerror* method which uses the error handler to collect errors reported. The arguments to the *AddError* method includes an error message generated by the parser and an *MPLEX.Parser.LexSpan* object representing the location of the error as well as the input stream buffer in which the error token is located.

The errors collected are available through a collection property on the *Parser* facade class, as described in Section 4.7.

## 4.6   Handling Escape Sequences

The chapters A.2 and A.7 of the SPARQL specification[1] describes codepoint escape sequences and common escape sequences respectively. An escape sequence is merely a substitute representation of an actual character. In SPARQL, all escape sequences start with a backslash (\) followed by en escape specifier.

The processes of handling escape sequences are described in the next sections and the invocation of these processes are made by the *Parser* facade class described in Section 4.7.

### 4.6.1   Codepoint Escape Sequences

The Unicode standard specifies unique identifiers for an inconceivable amount of characters and symbols. These positive integer identifiers are

| Escape | Codepoint |
|---|---|
| \u**xxxx** | A Unicode codepoint in the range $[0,\text{FFFF}_{16}]$ corresponding to the hexadecimal value (**xxxx**) |
| \U**xxxxxxxx** | A Unicode codepoint in the range $[0,\text{10FFFF}_{16}]$ corresponding to the hexadecimal value (**xxxxxxx**) |

Table 4.1: Codepoint escape sequence formats supported by SPARQL [1]

```
<ab\u00E9xy> # 00E9₁₆ is Latin small e with acute – <abéxy>
\u03B1:a     # 03B1₁₆ is Greek small alpha – α:a
a\u003Ab     # 003A₁₆ is colon – a:b
```

Figure 4.30: Usage examples of codepoint escape sequences

commonly referred to as *Unicode codepoints* [22]. Anywhere in a SPARQL query, a codepoint escape sequence may occur as a substitute for the character represented by that codepoint. Such an escape sequence can take two forms, as shown in Table 4.1. Figure 4.30 shows examples of codepoint escape sequence usage. Note that this implementation only supports codepoints in the range $[0,\text{FFFF}_{16}]$ in both forms, for the reasons explained in section 6.3.1.

Given that these codepoint escape sequences may occur anywhere in a SPARQL query, they should be processed and resolved before the query is being parsed. A preprocessor was developed to analyze a string input, replacing all occurrences of codepoint escape sequences with the actual character corresponding to the codepoint.

Two methods for resolving codepoint escape sequences were considered: regular expression (regex) transformation and manual analysis. While a regex transformation would result in less code, the regex pattern would be rather complex and result in poor performance compared to manual string analysis. Thus, the latter method was chosen. The algorithm used is roughly summarized in Figure 4.31.

The actual implementation was made in the method *ResolveCodepoint-EscapeSequences* of the class *SharQL.Utils.EscapeSequenceResolver*. *Argument-Exception*s are thrown if the codepoint is not recognized. Escape sequences that do not start with \u or \U are ignored and left for later processing, described next.

### 4.6.2 Common Escape Sequences

Common escape sequences are shorthand notations for the most common codepoint escape sequences and may only occur within string literals of a SPARQL query. Eight different common escape sequences are available, as

```
1. Create an empty string S.
2. For each character in input string I:
   2.1. If not a backslash, append character to S.
   2.2. If backslash, read next character:
        2.2.1 If lowercase 'u', read next 4 characters,
              resolve codepoint and append to S.
        2.2.2 If uppercase 'U', read next 8 characters,
              resolve codepoint and append to S.
        2.2.3 Else, append backslash and character to S
              for later processing.
3. Return S.
```

Figure 4.31: Algorithm for processing codepoint escape sequences

| Escape | Corresponding codepoint escape |
|--------|--------------------------------|
| \t | \u0009 (tab) |
| \n | \u000A (line feed) |
| \r | \u000D (carriage return) |
| \b | \u0008 (backspace) |
| \f | \u000C (form feed) |
| \" | \u0022 (quotation mark, double quote mark) |
| \' | \u0027 (apostrophe-quote, single quote mark) |
| \\ | \u005c (backslash) |

Table 4.2: Common escape sequences supported by SPARQL [1]

shown in Table 4.2.

Given that these escape sequences may only occur within string literals, they are most easily processed after parsing is completed. An algorithm similar to the one from Figure 4.31 was implemented in the method *Resolve-CommonEscapeSequences* of the class *SharQL.Utils.EscapeSequenceResolver*.

To traverse the resulting AST and resolve all common escape sequences in string literals, a visitor class was developed. The class *SharQL.Ast.-Visitor.EscapeSequenceResolverVisitor* implements the interface *IParserVisitor* and applies the aforementioned method for resolving common escape sequences to all string literals. An extract of the class is shown in Figure 4.32.

## 4.7 The Parser Facade Class

The MPLex and MPPG tools generate a scanner and a parser, respectively. In order for the end-user to perform any parsing, a scanner and a parser have to be instantiated and the parser has to be aware of the existence of the scanner instance in order to read its output. The scanner, in turn, needs to know what input to tokenize.

34

```
public Action<INode> this[string nodeType]
{
  get
  {
    if (nodeType == new StringLiteralNode().NodeType)
      return resolveEscapeSequence;
    else
      return noOp;
  }
}

private void resolveEscapeSequence(INode node)
{
  StringLiteralNode n = node as StringLiteralNode;
  if (n != null)
  {
    string escapedString = n.Value;
    string resolvedString = Utils.EscapeSequenceResolver.
        ResolveCommonEscapeSequences(escapedString);
    n.Value = resolvedString;
  }
}

private void noOp(INode node)
{
  // Ignore
}
```

Figure 4.32: Extract from *SharQL.Ast.Visitor.EscapeSequenceResolverVisitor*

```
┌─────────┐
│ SharQL  │
┌┴─────────┴──────────────────────────┐
│                                      │
│   ┌──────────────────────────────┐   │
│   │           Parser             │   │
│   ├──────────────────────────────┤   │
│   │ + Errors:Collection<Error> {get;} │   │
│   │ + AstRoot:INode {get;}       │   │
│   ├──────────────────────────────┤   │
│   │ + Parser()                   │   │
│   │ + Parse(source:string):bool  │   │
│   │ + Parse(source:string, offset:int):bool │   │
│   └──────────────────────────────┘   │
│                                      │
└──────────────────────────────────────┘
```

Figure 4.33: The *Parser* facade class

To encapsulate all this plumbing, a parser facade class, as shown in Figure 4.33, was developed. The intention of this class is to let the end-user instantiate it once, and make successive calls to one of its *Parse* methods without needing to have a notion of a scanner at all. When a parsing succeeds, the resulting AST is made available through the *AstRoot* property. Additionally, the parser facade reports any errors that may occur, through its *Errors* property.

Naming the facade class *Parser* for clarity was highly desirable. However, this crashed with the name that MPPG by default assigns its generated parser. Our solution to this was to modify the parser specification file by including a *%parsertype* specifier, thus renaming the resulting type of the parser class.

A UML sequence diagram showing parser facade class operation during parsing is shown in Figure 4.34.

## 4.8 Creating the SharQL Test Client

The intention of the SharQL Test Client is to have easy access to the SharQL parser library at any given time, and to be able to visualize the abstract syntax trees produced by the parser.

Given that SharQL is a .NET library, *Windows Forms* was chosen as platform for the test client. Visual Studio has great Windows Forms support, and creating graphical user interfaces is easily accomplished. So, a new Windows Forms project was created inside the SharQL solution, and the references required by SharQL where added, as shown in Figure 4.35.

The user interface was then designed as shown in Figure 4.36. The idea is that, once the query is parsed, the AST is displayed in the *TreeView* control to the left. When one of the nodes in the tree view is selected, details for

User

:Parser

new

:ParsingEngine

new

Parse(source:string)

:Scanner

new

:SharQLError-Handler

new

Escape-Sequence-Resolver

ResolveCodepointEscapeSequences(source:string):string

unescapedSource:string

SetSource(unescapedSource:string)

Handler = :SharQLErrorHandler

scanner = :Scanner

Performed via a visitor object that only handles string literal nodes.

Parse()

success

ResolveCommonEscapeSequences()

true

error

AddError(:Error)

false

Figure 4.34: UML sequence diagram showing parser facade class operation

Figure 4.35: Necessary references for the test client

that node is displayed in the text box to the right. Any errors reported by the parser is shown in the lower pane.

The test client can also perform a scan-only operation on the input file. By clicking the *Scan* button instead of the *Parse* button, the input is scanned and the stream of tokens produced by the scanner is displayed in a separate window.

Visualizing the abstract syntax tree is a matter of translating it to a tree of *TreeNode* objects. Once such a tree is created, it can be assigned to the *Nodes* property of the tree view in the test client.

Translating a node from the AST to a TreeNode object was done using the Visitor pattern explained in Section 4.4. Using this pattern, the translation was very simple. The *TreeBuilderVisitor* class first creates dictionary mapping INode objects to TreeNode objects. Remember that the visitor is traversing the tree depth-first. This means that when the visitor is applied to a node, it is safe to assume that the visitor has already been applied to all its child nodes.

Our visitor will not distinguish between different types of nodes. Thus, the same delegate is returned through the visitor indexer, regardless of the node type string provided. The implementation of the method represented by this delegate is shown in Figure 4.37. The TreeNode constructor accepts a textual caption as parameter. In this case, the class name is obtained from the NodeType property, assuming its default behavior of returning the full name of the class.

By assigning a reference to the actual AST node to the *Tag* property of the corresponding TreeNode object, detailed information about the node

38

Figure 4.36: The test client user interface

```
private void constructTreeNode(INode node)
{
  NodeBase nb = node as NodeBase;
  if (nb != null)
  {
    // TreeNodes maps INode objects to TreeNode objects
    TreeNodes[node] = new TreeNode(nb.NodeType.Split('.').Last());
    TreeNodes[node].Tag = node;
    foreach (INode child in nb.Children)
    {
      if (child != null && TreeNodes.ContainsKey(child))
      {
        TreeNodes[node].Nodes.Add(TreeNodes[child]);
      }
    }
  }
}
```

Figure 4.37: Code snippet from *TreeBuilderVisitor*

39

may be obtained later. Such information is needed when a node is selected. Then, a detailed description of the node is created and displayed in the right half of the user interface.

# Chapter 5

# Testing

The notion of testing is important in all software development projects. After all, without doing any kind of testing, one has no means of assessing whether the software works as it should. There are numerous ways of testing software. Black box testing is one way of performing tests, where the system as a whole is tested and its behavior compared to a reference of how it *should* behave. Another popular method is unit testing.

## 5.1   Unit Testing

Unit testing is a method of individually testing small parts, or units, of a larger system. Unit tests are usually based on a specification of how these units should behave. Often, the tests are even written before the work on implementing the unit has even begun. This is possible by defining interfaces that the units have to implement, and write tests against those interfaces.

It is important to write tests that cover as much as possible of the specification of how the unit should behave. After all, the intention of testing the behavior of the small units is to rule out subtle errors originating from a faulty part of a bigger system. By assuring that the individual units behave according to the specification, debugging integration errors in a later phase becomes more comprehensible.

Sometimes, the units are easily identified in the system specification. Otherwise, these units are usually identified and specified in the design phase. The SPARQL specification from W3C specifies the language, but does not consider the implementation of parsers. In this project, *MPLex* and *MPPG* generates a scanner and a parser, respectively. Thus, it is natural to consider the scanner as one unit and the parser as another unit. Testing the parser is a matter of feeding it with queries and to see that it parses those that are valid and refuses to parse queries that contain errors. Testing the scanner is a matter of verifying that it lexically separates queries into the

```
using Microsoft.VisualStudio.TestTools.UnitTesting;

(...)

[TestClass]
public class MyUnitTests
{
  (...)
}
```

Figure 5.1: Creating a test class

```
[TestClass]
public void MyUnitTest1()
{
  MyClass myObj = new MyClass();
  myObj.DoSomething(1, 2, 3);
  Assert.IsTrue(myObj.SomeProperty == 5);
  Assert.IsFalse(myObj.AnotherProperty > 0);
}
```

Figure 5.2: Creating a test method

appropriate tokens in the correct order.

## 5.2 Automated Testing in Visual Studio 2008

Visual Studio 2008 has integrated support for automated testing. By adding a *Test Project* to a Visual Studio solution, one can write tests in one of several programming languages and benefit from automatic execution and reporting based on the results.

To create unit tests in a test project, start by creating a public class and tag it with the *TestClass* attribute, as shown in Figure 5.1. The using statement shown resolves the namespace where the TestClass attribute resides. When Visual Studio is instructed to run the tests of the test projects, it uses reflection to decide which of the classes that are test classes.

Within a test class, an arbitrary number of test methods can be defined. These are methods that perform specific tasks and assert that the states and output of the tested objects are as expected. Similar to test classes, test methods require a *TestMethod* attribute. Assertions are performed using one of several static methods of the *Assert* class. An example test method is shown in Figure 5.2. This example creates an object and performs a specific operation to alter its state. Successively, the test performs two assertions on its state, according to some specification.

The static *Assert* class provides test authors with several assertions meth-

| Method name | Description |
|---|---|
| Are[Not]Equal | Verifies whether or not two objects are equal. |
| Are[Not]Same | Verifies whether or not two object variables refer to the same object. |
| Is[Not]InstanceOfType | Verifies whether or not an object is an instance of a specific type. |
| Is[Not]Null | Verifies whether or not an object is *null*. |
| IsFalse | Verifies that a condition is false. |
| IsTrue | Verifies that a condition is true. |
| Fail | Unconditionally fails the test. |
| Inconclusive | Indicates that an assertions can not be proved either true or false. |

Table 5.1: Methods available in *Assert* class

```
SELECT *
WHERE { }
```

Figure 5.3: The *syntax-basic-01* test from W3C

ods. These methods are summarized in Table 5.1. Most of the assertion methods are overloaded in order to support various data types and/or to let test authors provide optional error messages for the test reports. Where the test names in the table contains a "[Not]" part, a pair of methods exist, one containing "Not" and one that does not.

Visual Studio offers several ways of organizing tests and the execution of tests. For small projects like this one, however, it is typically desirable to run all the tests. This is done by clicking *Test → Run → All Tests in Solution*.

## 5.3 W3C SPARQL Test Suite

The *Data Access Working Group* of W3C has released a set of test cases, including syntax tests for the SPARQL language [2]. All in all, the test suite contains 199 SPARQL queries which should either parse or fail to parse. These tests are listed in Appendix E. While such tests do not fully test all aspects of a SPARQL parser, they constitute a solid foundation.

The tests vary from simple select queries as shown in Figure 5.3, to more complex queries as the one shown in Figure 5.4. Invalid queries, as in Figure 5.5, are also provided. Together, these queries cover most of the syntax that a SPARQL parser could encounter.

The test queries do not specify how a resulting AST should look like. After all, the structure of an AST depends on the application. Thus, it is

```
PREFIX rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>
CONSTRUCT { [] rdf:subject ?s ;
               rdf:predicate ?p ;
               rdf:object ?o }
WHERE {?s ?p ?o}
```

Figure 5.4: The *syntax-form-construct03* test from W3C

```
SELECT * WHERE {[] . }
```

Figure 5.5: The *syn-bad-bnode-dot* test from W3C

not possible to verify the correctness of the ASTs produced by our parser, based on these tests alone. The tests neither provide any means of testing the scanner isolatedly. Such tests have to be authored by manually identifying the token stream that should be produced. Once again, this is a natural constraint, given that the tokens chosen, like the AST, depends on the application.

### 5.3.1 Omission of Semantic Tests

Some of the syntax tests from the test suite actually test semantic matters. Table 5.2 lists the tests in question. All the tests are negative tests supposed to fail, but due to the current state of the parser, they all succeed. The reason why they should fail, is that the same blank node label is used across basic graph patterns, as shown in Figure 5.6. The parser is not supposed to care about semantics in this phase, and is for that reason handling the tests in questions as it should.

The fact that the tests listed in Table 5.2 are in fact semantic tests as well as syntactic tests, makes them unsuitable for testing the parser. For this reason, these tests have been omitted from the set of tests run when testing the parser. Of the total of 199 tests supplied, this leaves 188 tests that are still performed.

```
PREFIX : <http://xmlns.com/foaf/0.1/>

ASK { _:who :homepage ?homepage
      GRAPH ?g { ?someone :made ?homepage }
      _:who :schoolHomepage ?schoolPage }
```

Figure 5.6: The *syn-blabel-cross-graph-bad* test from W3C, semantically incorrect due to the use of the blank node *who* across basic graph patterns.

| Test File | Reason for failing |
|---|---|
| syntax-sparql3\syn-blabel-cross-graph-bad.rq<br>syntax-sparql3\syn-blabel-cross-optional-bad.rq<br>syntax-sparql3\syn-blabel-cross-union-bad.rq | Blank node *who* reused across basic graph patterns |
| syntax-sparql4\syn-bad-34.rq<br>syntax-sparql4\syn-bad-35.rq<br>syntax-sparql4\syn-bad-36.rq<br>syntax-sparql4\syn-bad-37.rq<br>syntax-sparql4\syn-bad-38.rq<br>syntax-sparql4\syn-bad-graph-breaks-BGP.rq<br>syntax-sparql4\syn-bad-opt-breaks-BGP.rq<br>syntax-sparql4\syn-bad-union-breaks-BGP.rq | Blank node *a* reused across basic graph patterns |

Table 5.2: Failing Semantic Tests Omitted From the W3C SPARQL Test Suite [2].

## 5.4   Automating the Test Suite

Using the integrated support for automated testing in Visual Studio 2008, the W3C SPARQL Test Suite may be completely automated. The running of the test suite is essential in determining the parser's level of conformance with the SPARQL grammar specification, and because of the great importance of running such tests, this should be as easy to carry out as possible.

The goal of automating the test suite is to be able to run all the tests in a single operation, and to receive a report stating the result of each and every test. The test suite contains 188 syntax tests relevant for the parser, each of which must be implemented by a corresponding test method as shown in Figure 5.2. The tests only differ in which syntax test file to parse. Thus, a great amount common code was identified, extracted and generalized to build a small framework for automating this specific test suite.

Using reflection features available in the C# programming language, a method is able to determine the name of the caller method. Figure 5.7 shows the code snippet necessary for accessing the caller stack and fetching the second top-most stack frame to extract the caller method name. By naming each test method according to the corresponding test syntax file and having each test method call a common method, this latter method may determine which test syntax file to load and can carry out the test on its own.

Naming the test methods according to the file names of the test syntax files supplied in the W3C SPARQL Test Suite is advantageous since the corresponding test syntax file for a failing test can very quickly be determined. In addition, since the test syntax file may be derived from the test method's name, the code necessary for calling the common test method from each

```
StackFrame stackFrame = new StackFrame(1, true);
MethodBase methodBase = stackFrame.GetMethod();
string methodName = methodBase.Name;
```

Figure 5.7: C# Caller Method Name Extraction

```
[TestMethod]
public void sparql1__basic_01() { Assert.IsTrue(
  ParseQueryFile("syntax")); }
```

Figure 5.8: Test Method Sample

test method may be reduced to a minimum. Figure 5.8 shows an actual test method calling the common *ParseQueryFile* method. The test method name is actually a composition of the test set name followed by the specific test name. Thus, *sparql1__basic_01* corresponds to the *basic_01* syntax test in the *sparql1* test set.

With syntax tests mapped to test methods as described, the entire test suite containing the 188 syntax tests distributed over five tests sets, may be performed by a single operation.

## 5.5   Custom Tests

In addition to the tests in the W3C SPARQL Test Suite [2], a set of custom tests have been implemented to test other aspects of the parser. The tests supplied in the test suite test the parser as a whole without drilling down on a specific component. The purpose of the custom tests are to test components independently, allowing for more fine grained testing.

### 5.5.1   Scanner Tests

As the purpose of the scanner is to identify tokens in the input stream, a natural approach for testing the scanner separately is to assure that the correct tokens are identified in the correct order. The test suite tests are roughly split into categories, each of which focus on a certain aspect of the SPARQL grammar. For each such category a custom test has been implemented that feeds a specific test to the scanner and validates the tokens returned. Both token types and positions should be correct in order for the test to succeed.

Figure 5.9 shows the custom scanner test for the SPARQL query shown in Figure 5.3. The tokens returned should be *SELECT*, *\**, *WHERE*, *{* and *}*, in that order.

```
[TestMethod]
public void Scanner_1_basic_01()
{
  Assert.IsTrue(assertScannerOutput(
    @"..\..\..\Tests\syntax-sparql1\syntax-basic-01.rq",
    new Tokens[]
    {
      Tokens.SELECT,
      (Tokens)((int)'*'),
      Tokens.WHERE,
      (Tokens)((int)'{'),
      (Tokens)((int)'}')
    })
  );
}
```

Figure 5.9: Custom Scanner Test

```
string escapedString = @"\u0061";
string expected = "a";
string actual;
actual = EscapeSequenceResolver.ResolveCodepointEscapeSequences(
  escapedString);
Assert.AreEqual(expected, actual);
```

Figure 5.10: Common escape sequence resolver test

### 5.5.2 Escape Tests

As described in Section 4.6, escape sequences are processed and resolved both ahead of and after the parsing. A set of custom tests have been implemented for testing this functionality separately as well. The escape sequence resolvers and the escape visitor are all tested separately.

Figure 5.10 shows an extract from the test for the common escape sequence resolver test

# Chapter 6

# Results

This chapter starts by illustrating, step by step, how the resulting SharQL parser handles a sample SPARQL query. Next, the results from running the SharQL parser against the *W3C SPARQL Test Suite* are presented. Finally, known nonconformities with the W3C SPARQL specification are explained.

## 6.1   A Simple Example

This section presents a complete example involving all essential steps in the parse process, beginning with the input SPARQL query and resulting in the AST generated by the parser.

Most implementation-specific details have been left out to as these are described in Chapter 4. The purpose of this example is to present a higher-level view of the parser and its usage.

### 6.1.1   Example Query

The example SPARQL query used is a rather simple query, as shown in Figure 6.1. The complexity of the query is reasonable in order to highlight the essential parsing steps at the desired level of detail.

```
PREFIX dc: <http://purl.org/dc/elements/1.1/>

SELECT ?title
WHERE
{
  <http://en.wikipedia.org/wiki/Tony_Benn> dc:title ?title
}
```

Figure 6.1: SPARQL example query, identical to the query presented in Figure 2.4

```
PREFIX
PNAME_NS
IRI_REF
SELECT
VAR1
'{'
IRI_REF
PNAME_LN
VAR1
'}'
```

Figure 6.2: SPARQL example query scanner generated token stream

### 6.1.2 The Scanner Token Generation

Given the example query in Figure 6.1 the scanner generates the token
stream shown in Figure 6.2. The tokens are passed to the parser one by
one in the order of appearance. This process is described in Section 4.3.2.

### 6.1.3 The AST generation

As the parser receives the tokens from the scanner, it collects tokens until a
rule in the parser specification is matched by the tokens available. When a
rule is matched, a corresponding AST class is instantiated to represent the
rule. The AST class hierarchy is described in Section 4.2.

In this example the *PREFIX* token is the first one collected. As this does
not match any rule by itself, the parser simply keeps hold of it while fetch-
ing the next token from the scanner.

The second token identified by the scanner is the *PNAME_NS* token,
which has a corresponding rule in the parser specification, shown in Figure
6.3. When the rule is matched, a *PrefixDeclNode* object is created represent-
ing the token. The *sc.yytext* expression represents the textual representation
of the token, being the namespace name "'dc'".

The third token is the IRI_REF token. The rule matched for this token is
shown in Figure 6.4. The corresponding object created for this token is an
*IriRefNode* object.

Both the objects created replace their corresponding tokens in the parser
collection of tokens received from the scanner, both represented as the match-
ing rule from which they were constructed. At this point, the collection con-
tains a *PREFIX* token, a *PrefixDeclNode* represented by the *Pname* rule and
an *IriRefNode* represented by the *IriRef* rule. Those three elements match
the *Prefix* rule, shown in Figure 6.5.

The *Prefix* rule is handed the elements in the parser collection and uses
them to compose its corresponding object. The remaining tokens are han-
dled in a similar manner, until the entire token stream has been processed

```
Pname
: PNAME_NS
{
  $$ = new PrefixDeclNode(sc.yytext, null);
}
;
```

Figure 6.3: *PNAME_NS* Token Parser Rule

```
IriRef
: IRI_REF
{
  $$ = new IriRefNode() { Value = sc.yytext };
}
;
```

Figure 6.4: *IRI_REF* Token Parser Rule

resulting in one *QueryNode* object representing the entire query. This node is the root of the generated AST, as shown in Figure 6.6.

The relationships between the entities involved in the parsing of a SPARQL query is shown in Figure 6.7.

## 6.2 Test Results

Testing has been an important tool in the development of the SharQL parser, both as a measure of correctness and as a bug tracking mechanism. The W3C SPARQL Test Suite [2] has been a great resource for validating the correctness of the parser as a whole through black box testing. The tests from W3C are listed in Appendix E. Custom tests have helped testing single components through unit testing.

```
/* [4] */
PrefixDecl
: PREFIX Pname IriRef
{
  PrefixDeclNode pn = (PrefixDeclNode)$2;
  $$ = new PrefixDeclNode(pn.Namespace, $3);
}
;
```

Figure 6.5: *Prefix* Parser Rule

51

Figure 6.6: Example generated AST

Figure 6.7: Relationships between the parser context entities

| | Test Set 1 | Test Set 2 | Test Set 3 | Test Set 4 | Test Set 5 |
|---|---|---|---|---|---|
| ■ Failed | 0 | 0 | 3 | 8 | 0 |
| ■ Succeeded | 81 | 53 | 48 | 4 | 2 |

Figure 6.8: Test Suite Results

### 6.2.1  W3C SPARQL Test Suite

The tests described in Section 5.3 were run against our parser, resulting in a parse success or parse failure. Test queries that should parse, but did not, indicated errors in our parser. During development, quite a few errors were located by means of valid queries which did not parse successfully.

On the occurrence of a test failing, the test query in question has been examined and used as a basis for tracking the creation of the corresponding AST during parsing to help identifying the point of failure. In most cases this has lead to an error being identified and corrected.

A small amount of tests were identified as semantical tests as well as syntactical tests, despite being classified as syntax tests in the test suite. However, the parser does currently not operate on a semantical level, and the tests in question have thus been excluded, as discussed in Section 5.3.1.

The test suite is divided into five test sets, each containing a number of single test queries. Figure 6.8 shows the results of running all the tests supplied in the test suite, grouped by test set. The failing tests represent the semantical tests that have now been excluded from the test setup.

## 6.3 Specification Nonconformities

### 6.3.1 Unsupported Unicode Codepoints

A SPARQL query string is a Unicode character string [1]. As discussed in Section 4.6.1, the Unicode standard specifies unique identifiers (called codepoints) for all Unicode characters. The set of defined codepoints is divided into different *planes*. Each plane consists of 65,536 ($2^{16}$) codepoints. Plane 0 is the range $[0000_{16}, \text{FFFF}_{16}]$, Plane 1 is $[10000_{16}, \text{1FFFF}_{16}]$ etc. Plane 0 is called the *Basic Multilingual Plane* (BMP), while codepoints from the remaining planes are called *Supplementary Code Points*. [23]

The basic unit recognized by the MPLex generated scanner is the .NET *System.Char* type. This type represents a 16 bit integer value. Thus, a single *System.Char* object can only represent a Unicode character in the Basic Multilingual Plane. To represent Unicode characters from the remaining planes, two successive *System.Char* objects are required. This leads to undefined behavior by the MPLex generated scanner.

While a SPARQL query, according to W3C's specification, may contain characters from all Unicode planes, the SharQL parser developed in this project only has a defined behavior for queries containing characters from the Basic Multilingual Plane. While this violates the SPARQL specification, we consider it a minor issue, as the majority of common-use characters fit into the BMP [22].

### 6.3.2 Semantic Specifications

According to the SPARQL specification, Section A.6: *"The same blank node label may not be used in two separate basic graph patterns with a single query"* [1].

As discussed in Section 5.3.1, this requirement is considered a semantic rule and thus outside the scope of this project. The SharQL parser will not yield any errors when parsing SPARQL queries reusing the same blank node label in two separate basic graph patterns.

# Chapter 7

# Conclusion and Further Work

This chapter briefly discusses the course of the project in terms of key ingredients and decisions, resulting in the final outcome. Also, the planned further work is presented, placing the project in its superior context.

## 7.1 Conclusion

The project has been greatly influenced by the tools used, being Visual Studio and the Managed Babel package containing the MPLex and MPPG lexer and parser generator. Allowing for the lexer and parser generator to be integrated in the SharQL project itself has greatly simplified the parser generation by, in reality, making it completely transparent.

As the SPARQL grammar is defined on *Extended Backus-Naur Form* (EBNF), it had to be translated into the BNF-based specification languages used by MPLex and MPPG. This has introduced quite a few additional grammar rules to compensate for the lack of expressive power in the specification languages. Besides being the most time-consuming task in the project altogether, this translation has been the number one source of bugs during the development.

Conformance with the SPARQL grammar is essential, as well as keeping the implementation free of bugs in general. The W3C SPARQL Test Suite has been an important tool in the strive for correct behaviour. The suite contains nearly 200 tests relevant for the SharQL parser, testing different aspects of the grammar. However, despite being classified as syntactic tests, a rather small amount of the tests actually tested semantic matters as well. As such matters are not taken into account at this stage, the tests in question have been left out from the set of tests used to validate the conformity of the parser.

The outcome of the project is the SharQL parser. It runs on the Microsoft .NET framwork, and generates abstract syntax trees from SPARQL queries. The only nonconformity with the SPARQL specification is the lack of sup-

port for Unicode characters that exceed 16 bits in size. All relevant tests from the W3C SPARQL Test Suite pass without exceptions, indicating the level of conformity with the SPARQL specification.

## 7.2 Further Work

This project and its outcome directly prepares for a master thesis for the iAD Research Centre in spring 2009. The thesis will be part of the MARS project, and addresses the following two problems:

- Efficiently transforming RDF data to, and representing it in, the representation used in MARS engine.

- Generating an executable query plan from a SPARQL query for the MARS .NET engine.

RDF is verbose by nature, especially in its XML-based representation. Regardless of the representation used, triples still mainly consist of URIs, causing a great deal of redundancy. Also, RDF data represents a graph of nodes, which is quite different from typical relational data representations. These properties may call for a different approach in storing and indexing the data in order to allow for queries to be executed efficiently.

When a scheme for representing the RDF-data has been defined, the next step is to generate executable query plans. The outcome of this project is the SharQL parser currently generating abstract syntax trees from SPARQL queries. The ASTs serve as a basis for such query plans, and should be transformed from the current tree of nodes to a tree of query elements defined by the MARS .NET project.

# References

[1] W3C. SPARQL query language for RDF. `http://www.w3.org/TR/2008/REC-rdf-sparql-query-20080115/`, Jan. 2008.

[2] W3C. DAWG Testcases. `http://www.w3.org/2001/sw/DataAccess/tests/r2`.

[3] W3C. W3C semantic web activity. `http://www.w3.org/2001/sw/`, July 2008.

[4] W3C. RDF primer. `http://www.w3.org/TR/2004/REC-rdf-primer-20040210/`, Feb. 2004.

[5] DCMI-Libraries Working Group. Library application profile. `http://dublincore.org/documents/2004/09/10/library-application-profile/`, Sept. 2004.

[6] Turtle - Terse RDF Triple Language. `http://www.w3.org/TeamSubmission/2008/SUBM-turtle-20080114/`, Jan. 2008.

[7] Notation 3. `http://www.w3.org/DesignIssues/Notation3.html`, Mar. 2006.

[8] W3C. RDF Data Access Use Cases and Requirements. `http://www.w3.org/TR/2005/WD-rdf-dawg-uc-20050325/`, Mar. 2005.

[9] W3C. W3C Opens Data on the Web with SPARQL. `http://www.w3.org/2007/12/sparql-pressrelease`, January 2008.

[10] W3C. Testimonials for "W3C Opens Data on the Web with SPARQL" Press Release. `http://www.w3.org/2007/12/sparql-testimonial`, 2008.

[11] Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman. *Compilers: Principles, Techniques, and Tools*. Addison Wesley, 2nd edition, 2006.

[12] ANother Tool for Language Recognition (ANTLR). `http://www.antlr.org/`, Oct. 2009.

[13] Compiler compiler/Recursive descent (Coco/R). `http://www.ssw.uni-linz.ac.at/Research/Projects/Coco/`, Oct. 2009.

[14] Microsoft. Managed babel. `http://msdn.microsoft.com/en-us/library/bb165037(VS.90).aspx`, Nov. 2007.

[15] Queensland University of Technology. The Gardens Point Scanner Generator (GPLEX). `http://plas.fit.qut.edu.au/gplex/`, Jan. 2007.

[16] Queensland University of Technology. The Gardens Point Parser Generator (GPPG). `http://plas.fit.qut.edu.au/gppg/`, Jan. 2007.

[17] The Lex & Yacc page. `http://dinosaur.compilertools.net/`.

[18] Microsoft. Structs (C#). `http://msdn.microsoft.com/en-us/library/saxz13w4(VS.85).aspx`, Nov. 2008.

[19] Microsoft. MSDN Academic Alliance. `http://msdn.microsoft.com/en-us/academic/`, Nov. 2008.

[20] W3C. Extensible Markup Language (XML) 1.1. `http://www.w3.org/TR/2004/REC-xml11-20040204/#sec-notation`, February 2004.

[21] Erich Gamma, Richard Helm, Ralph Johnson, and John M. Vlissides. *Design Patterns: Elements of Reusable Object-Oriented Software*. Addison-Wesley Professional, 1994.

[22] The Unicode Consortium. About the Unicode Standard. `http://www.unicode.org/standard/standard.html`, Nov. 2008.

[23] The Unicode Consortium. Glossary of Unicode Terms. `http://unicode.org/glossary/`, Nov. 2008.

# Glossary

**AST** Abstract Syntax Tree, a tree representation of the syntax of some source code.

**BNF** Backus-Naur Form, a metasyntax used to express context-free grammars.

**C#** One of the programming languages supported by Microsoft for creating applications that execute on the .NET Framework.

**EBNF** Extended Backus-Naur Form, an extension to BNF containing several shorthand notations.

**Escape Sequence** In the context of this project, an escape sequence is a way of representing Unicode characters by appending a character specifier to an escape character.

**Grammar Production** A grammar production specifies possible substitutions for a specific grammar symbol.

**iAD** Information Access Disruptions, a constellation between FAST, two Norwegian enterprises and different research environments with a goal to develop the best search technology in the world.

**IRI** Internationalized Resource Identifier, a generalization of the Uniform Resource Identifier (URI), allowing Unicode character rather than being restricted to a subset of ASCII characters.

**LALR Parser** Lookahead LR parser, produces the rightmost derivation, reading the input from left to right.

**LL(k) Parser** Produces the leftmost derivation, reading the input from left to right, using at most $k$ tokens lookahead.

**Managed Babel** A package for integrating new languages into Microsoft's Visual Studio software. Contains tools that may also be used for general parsing purposes.

**.NET Framework** A software technology available from Microsoft which provides a library of pre-coded components and a virtual machine for managing application execution.

**Parser** A program that performs syntactical analysis on a sequence of tokens to determine the grammatical structure. A parser usually produces an Abstract Syntax Tree (AST) for further analysis.

**RDF** Resource Description Framework, a model for representing information about resources on the World Wide Web.

**Scanner** A program that reads a sequence of characters and produce a sequence of tokens which represent one or more characters. A scanner performs the first step (the lexical analysis) when parsing input in a given language.

**Semantic Web** An extension of the World Wide Web in which the semantics of information and services on the web is defined, making it possible for the web to understand and satisfy the requests of people and machines to use the web content.

**SharQL** The SPARQL parser created in this project.

**SPARQL** SPARQL Protocol And RDF Query Language, a query language for RDF data.

**Token** A categorized block of text, recognized by a scanner. A scanner outputs an array of tokens which may be interpreted by a parser.

**Turtle** A serialization format for RDF, resulting in a less verbose output than the equivalent XML serialization.

**Unicode** A specification assigning unique code points to symbols from most of the world's writing systems.

**Unit Test** A test that verifies that an individual unit of source code is working properly.

**URI** Uniform Resource Identifier, a compact string of characters used to identify or name a resource on the Internet.

**Visitor Pattern** A way of separating an algorithm from an object structure upon which it operates.

**Visual Studio** An integrated development environment (IDE) from Microsoft with support for developing .NET Framework applications.

**W3C** World Wide Web Consortium, the main international standards organization for the World Wide Web.

**W3C SPARQL Test Suite**  A set of tests provided by SPARQL to help ensure that a SPARQL parser behaves according to the specification.

# Appendix A

# XML for Project File

```
1   <UsingTask TaskName="MPLex"
2       AssemblyFile="lib\Microsoft.VsSDK.Build.Tasks.dll" />
3   <UsingTask TaskName="FindVsSDKInstallation"
4       AssemblyFile="lib\Microsoft.VsSDK.Build.Tasks.dll" />
5   <!--Set the general properties for this installation of the SDK-->
6   <PropertyGroup>
7     <VsSDKVersion>9.0</VsSDKVersion>
8     <VSSDKTargetPlatformVersion>9.0</VSSDKTargetPlatformVersion>
9     <VSSDKTargetPlatformRegRoot>
10      Software\Microsoft\VisualStudio\$(VSSDKTargetPlatformVersion)
11    </VSSDKTargetPlatformRegRoot>
12  </PropertyGroup>
13  <Target Name="FindSDKInstallation"
14      Condition="'$(VsSDKInstall)'==''">
15    <FindVsSDKInstallation SDKVersion="$(VsSDKVersion)">
16      <Output TaskParameter="InstallationPath"
17          PropertyName="VsSDKInstall" />
18      <Output TaskParameter="IncludesPath"
19          PropertyName="VsSDKIncludes" />
20      <Output TaskParameter="ToolsPath"
21          PropertyName="VsSDKToolsPath" />
22    </FindVsSDKInstallation>
23  </Target>
24  <PropertyGroup>
25    <TargetVSVersion Condition="'$(TargetVSVersion)' == ''">
26      $(VSSDKTargetPlatformVersion)
27    </TargetVSVersion>
28  </PropertyGroup>
29  <!--
30    ========================================================
31                Generate code from LEX files
32    ========================================================
33  -->
34  <PropertyGroup>
35    <!--Make sure that the lexer runs before the C# compiler-->
36    <CoreCompileDependsOn>
37      $(CoreCompileDependsOn);GenerateCodeFromLex
38    </CoreCompileDependsOn>
39  </PropertyGroup>
40  <Target Name="GenerateCodeFromLex"
41      Condition="'$(BuildingProject)'!='false'"
42      Inputs="@(MPLexCompile);$(LexFrameFile)"
43      Outputs="@(MPLexCompile->'$(IntermediateOutputPath)%(FileName).cs')"
```

```
44         DependsOnTargets="$(GenerateCodeFromLexDependsOn)">
45      <MPLex InputFile="@(MPLexCompile)"
46         OutputFile="@(MPLexCompile->'$(IntermediateOutputPath)%(FileName).cs')"
47         CompressTables="$(CompressTables)" FrameFile="$(LexFrameFile)"
48         SDKVersion="$(VsSDKVersion)">
49        <Output TaskParameter="OutputFile" ItemName="Compile" />
50        <Output TaskParameter="OutputFile" ItemName="FileWrites" />
51      </MPLex>
52    </Target>
53    <!--
54      ============================================================
55                  Generate parser code from Y files
56      ============================================================
57    -->
58    <PropertyGroup>
59      <GenerateParserCodeFromGrammarDependsOn>
60        $(GenerateParserCodeFromGrammarDependsOn);FindSDKInstallation
61      </GenerateParserCodeFromGrammarDependsOn>
62      <CoreCompileDependsOn>
63        $(CoreCompileDependsOn);GenerateParserCodeFromGrammar
64      </CoreCompileDependsOn>
65    </PropertyGroup>
66    <Target Name="GenerateParserCodeFromGrammar"
67        Condition="'$(BuildingProject)'!='false'"
68        Inputs="@(MPPGCompile)"
69        Outputs="@(MPPGCompile->'$(IntermediateOutputPath)%(FileName).cs')"
70        DependsOnTargets="$(GenerateParserCodeFromGrammarDependsOn)">
71      <!--Check if there are .lex files in the project because in
72          this case mppg should generate the base classes
73          used by the code generated by MPLex-->
74      <CreateProperty Value="-mplex" Condition="'@(MPLexCompile)' != ''">
75        <Output TaskParameter="ValueSetByTask"
76            PropertyName="__GenerateForMPLex" />
77      </CreateProperty>
78      <!--Run the command line tool that generates the cs files.-->
79      <!--Exec Command attribute should not span several lines in actual
80          project file-->
81      <Exec Command="&quot;$(VsSDKToolsPath)\MPPG.exe&quot;
82   $(__GenerateForMPLex) @(MPPGCompile->'&quot;%(Identity)&quot;') &gt;
83   @(MPPGCompile->'&quot;$(IntermediateOutputPath)%(FileName).cs&quot;')" />
84      <!--Add the generated files to the list of the files to compile.-->
85      <CreateItem
86          Include="@(MPPGCompile->'$(IntermediateOutputPath)%(FileName).cs')">
87        <Output TaskParameter="Include" ItemName="Compile" />
88        <Output TaskParameter="Include" ItemName="FileWrites" />
89      </CreateItem>
90    </Target>
```

# Appendix B

# Scanner Specification

```
1    %using SharQL.Utils;
2
3    /* Defines the generated scanner namespace. */
4    %namespace SharQL
5
6    /* Enables unicode characters. */
7    %option unicode
8
9    %{
10
11     // Called each time a token has been returned.
12     internal void LoadYylval()
13     {
14       // Collects token location information and makes it available for the
15       // parser through the yylloc variable.
16       yylloc = new CustomLexLocation(yyline, yycol, yyline, yycol + yyleng);
17     }
18
19     // Called by the parser on the occurence of an unexpected or missing
20     // token.
21     public override void yyerror(string s, params object[] a)
22     {
23       // Adds the reported error to the error handler.
24       handler.AddError(s, new MPLEX.Parser.LexSpan(tokLin, tokCol, tokLin,
25           tokECol, tokPos, tokEPos, new ScanBuffProxy(buffer)));
26     }
27
28     // Proxy class for the MPLEX.Lexer.ScanBuff class to allow for the
29     // MPLEX.Parser.LexSpan class to be used with the generated
30     // SharQL.ScanBuff class.
31     public class ScanBuffProxy : MPLEX.Lexer.ScanBuff
32     {
33       private ScanBuff buffer;
34
35       public override int Pos
36       {
37         get { return buffer.Pos; }
38         set { buffer.Pos = value; }
39       }
40
41       public override int Read()
42       {
43         return buffer.Read();
```

67

```
43        }
44
45      public override int Peek()
46      {
47        return buffer.Peek();
48      }
49
50      public override int ReadPos
51      {
52        get { return buffer.ReadPos; }
53      }
54
55      public override string GetString(int b, int e)
56      {
57        return buffer.GetString(b, e);
58      }
59
60      public ScanBuffProxy(ScanBuff buffer)
61      {
62        this.buffer = buffer;
63      }
64    }
65
66  %}
67
68  /*
69   * SPARQL Query Language for RDF
70   * W3C Recommendation 15 January 2008
71   *
72   * http://www.w3.org/TR/2008/REC-rdf-sparql-query-20080115/
73   *
74   */
75
76  iri_ref               <([^<>"{}|^`\\\u0000-\u0020])*>
77  langtag               @[a-zA-Z]+(\-[a-zA-Z0-9]+)*
78  exponent              [eE][\+\-]?[0-9]+
79  integer               [0-9]+
80  integer_positive      \+{integer}
81  integer_negative      \-{integer}
82  decimal               [0-9]+\.[0-9]*|\.[0-9]+
83  decimal_positive      \+{decimal}
84  decimal_negative      \-{decimal}
85  double                [0-9]+\.[0-9]*{exponent}|
        \.([0-9])+{exponent}|([0-9])+{exponent}
86  double_positive       \+{double}
87  double_negative       \-{double}
88  echar                 \\[tbnrf\\"']

90  /* Range 10000-EFFFF is left out, as MPLEX currently only supports 16-bit
        unicode characters. */
91  pn_chars_base         [A-Z]|[a-z]|[\u00C0-\u00D6]|[\u00D8-\u00F6]|
        [\u00F8-\u02FF]|[\u0370-\u037D]|[\u037F-\u1FFF]|
        [\u200C-\u200D]|[\u2070-\u218F]|[\u2C00-\u2FEF]|
        [\u3001-\uD7FF]|[\uF900-\uFDCF]|[\uFDF0-\uFFFD]
92
93  pn_chars_u            {pn_chars_base}|_
94  pn_chars              {pn_chars_u}|\-|[0-9]|\u00B7|
        [\u0300-\u036F]|[\u203F-\u2040]
95  pn_prefix             {pn_chars_base}(({pn_chars}|\.)*{pn_chars})?
96  pn_local              ({pn_chars_u}|[0-9]) (({pn_chars}|\.)*{pn_chars})?
97  pname_ns              {pn_prefix}?:
98  pname_ln              {pname_ns}{pn_local}
```

```
 99  blank_node_label      _:{pn_local}
100  varname               ({pn_chars_u}|[0-9])({pn_chars_u}|[0-9]|
         \u00B7|[\u0300-\u036F]|[\u203F-\u2040])*
101  var1                  \?{varname}
102  var2                  \${varname}
103  string_literal1       '(([^\u0027\u005C\u000A\u000D])|{echar})*'
104  string_literal2       \"(([^\u0022\u005C\u000A\u000D])|{echar})*\"
105  string_literal_long1  '''((''|'')?([^'\\]|{echar}))*'''
106  string_literal_long2  \"\"\"((\"|\"\")?([^\"\\]|{echar}))*\"\"\"
107  ws                    \u0020|\u0009|\u000D|\u000A
108  nil                   \(({ws}*\)
109  anon                  \[{ws}*\]
110  eq                    "="
111  less                  "<"
112  greater               ">"
113  or                    "||"
114  and                   "&&"
115  not                   "!="
116  lessorequal           "<="
117  greaterorequal        ">="
118  hats                  "^^"
119  comment               #[^\n]*
120
121  /*
122   * A.8 Grammar
123   *
124   * (...) Keywords are matched in a case-insensitive manner with the
125   * exception of the keyword 'a'.
126   */
127  base                  [Bb][Aa][Ss][Ee]
128  prefix                [Pp][Rr][Ee][Ff][Ii][Xx]
129  select                [Ss][Ee][Ll][Ee][Cc][Tt]
130  distinct              [Dd][Ii][Ss][Tt][Ii][Nn][Cc][Tt]
131  reduced               [Rr][Ee][Dd][Uu][Cc][Ee][Dd]
132  construct             [Cc][Oo][Nn][Ss][Tt][Rr][Uu][Cc][Tt]
133  describe              [Dd][Ee][Ss][Cc][Rr][Ii][Bb][Ee]
134  ask                   [Aa][Ss][Kk]
135  from                  [Ff][Rr][Oo][Mm]
136  named                 [Nn][Aa][Mm][Ee][Dd]
137  where                 [Ww][Hh][Ee][Rr][Ee]
138  orderby               [Oo][Rr][Dd][Ee][R]\u0020[Bb][Yy]
139  asc                   [Aa][Ss][Cc]
140  desc                  [Dd][Ee][Ss][Cc]
141  limit                 [Ll][Ii][Mm][Ii][Tt]
142  offset                [Oo][Ff][Ff][Ss][Ee][Tt]
143  optional              [Oo][Pp][Tt][Ii][Oo][Nn][Aa][Ll]
144  graph                 [Gg][Rr][Aa][Pp][Hh]
145  filter                [Ff][Ii][Ll][Tt][Ee][Rr]
146  union                 [Uu][Nn][Ii][Oo][Nn]
147  str                   [Ss][Tt][Rr]
148  lang                  [Ll][Aa][Nn][Gg]
149  langmatches           [Ll][Aa][Nn][Gg][Mm][Aa][Tt][Cc][Hh][Ee][Ss]
150  datatype              [Dd][Aa][Tt][Aa][Tt][Yy][Pp][Ee]
151  bound                 [Bb][Oo][Uu][Nn][Dd]
152  sameterm              [Ss][Aa][Mm][Ee][Tt][Ee][Rr][Mm]
153  isiri                 [Ii][Ss][Ii][Rr][Ii]
154  isuri                 [Ii][Ss][Uu][Rr][Ii]
155  isblank               [Ii][Ss][Bb][Ll][Aa][Nn][Kk]
156  isliteral             [Ii][Ss][Ll][Ii][Tt][Ee][Rr][Aa][Ll]
157  regex                 [Rr][Ee][Gg][Ee][Xx]
158  true                  [Tt][Rr][Uu][Ee]
159  false                 [Ff][Aa][Ll][Ss][Ee]
```

```
160
161   %%
162
163   {base}                 {return (int)Tokens.BASE;}
164   {prefix}               {return (int)Tokens.PREFIX;}
165   {select}               {return (int)Tokens.SELECT;}
166   {distinct}             {return (int)Tokens.DISTINCT;}
167   {reduced}              {return (int)Tokens.REDUCED;}
168   {construct}            {return (int)Tokens.CONSTRUCT;}
169   {describe}             {return (int)Tokens.DESCRIBE;}
170   {ask}                  {return (int)Tokens.ASK;}
171   {from}                 {return (int)Tokens.FROM;}
172   {named}                {return (int)Tokens.NAMED;}
173   {where}                {return (int)Tokens.WHERE;}
174   {orderby}              {return (int)Tokens.ORDERBY;}
175   {asc}                  {return (int)Tokens.ASC;}
176   {desc}                 {return (int)Tokens.DESC;}
177   {limit}                {return (int)Tokens.LIMIT;}
178   {offset}               {return (int)Tokens.OFFSET;}
179   {optional}             {return (int)Tokens.OPTIONAL;}
180   {graph}                {return (int)Tokens.GRAPH;}
181   {filter}               {return (int)Tokens.FILTER;}
182   {union}                {return (int)Tokens.UNION;}
183   {eq}                   {return (int)'=';}
184   {less}                 {return (int)'<';}
185   {greater}              {return (int)'>';}
186   {or}                   {return (int)Tokens.OR;}
187   {and}                  {return (int)Tokens.AND;}
188   {not}                  {return (int)Tokens.NOT;}
189   {lessorequal}          {return (int)Tokens.LESSOREQUAL;}
190   {greaterorequal}       {return (int)Tokens.GREATEROREQUAL;}
191   {str}                  {return (int)Tokens.STR;}
192   {lang}                 {return (int)Tokens.LANG;}
193   {langmatches}          {return (int)Tokens.LANGMATCHES;}
194   {datatype}             {return (int)Tokens.DATATYPE;}
195   {bound}                {return (int)Tokens.BOUND;}
196   {sameterm}             {return (int)Tokens.SAMETERM;}
197   {isiri}                {return (int)Tokens.ISIRI;}
198   {isuri}                {return (int)Tokens.ISURI;}
199   {isblank}              {return (int)Tokens.ISBLANK;}
200   {isliteral}            {return (int)Tokens.ISLITERAL;}
201   {regex}                {return (int)Tokens.REGEX;}
202   {true}                 {return (int)Tokens.TRUE;}
203   {false}                {return (int)Tokens.FALSE;}
204   {hats}                 {return (int)Tokens.HATS;}
205   {iri_ref}              {return (int)Tokens.IRI_REF;}
206   {langtag}              {return (int)Tokens.LANGTAG;}
207   {exponent}             {return (int)Tokens.EXPONENT;}
208   {integer}              {return (int)Tokens.INTEGER;}
209   {integer_positive}     {return (int)Tokens.INTEGER_POSITIVE;}
210   {integer_negative}     {return (int)Tokens.INTEGER_NEGATIVE;}
211   {decimal}              {return (int)Tokens.DECIMAL;}
212   {decimal_positive}     {return (int)Tokens.DECIMAL_POSITIVE;}
213   {decimal_negative}     {return (int)Tokens.DECIMAL_NEGATIVE;}
214   {double}               {return (int)Tokens.DOUBLE;}
215   {double_positive}      {return (int)Tokens.DOUBLE_POSITIVE;}
216   {double_negative}      {return (int)Tokens.DOUBLE_NEGATIVE;}
217   {echar}                {return (int)Tokens.ECHAR;}
218   {pname_ns}             {return (int)Tokens.PNAME_NS;}
219   {pname_ln}             {return (int)Tokens.PNAME_LN;}
220   {blank_node_label}     {return (int)Tokens.BLANK_NODE_LABEL;}
221   {var1}                 {return (int)Tokens.VAR1;}
```

70

```
222  {var2}                 {return (int)Tokens.VAR2;}
223  {string_literal1}      {return (int)Tokens.STRING_LITERAL1;}
224  {string_literal2}      {return (int)Tokens.STRING_LITERAL2;}
225  {string_literal_long1}{return (int)Tokens.STRING_LITERAL_LONG1;}
226  {string_literal_long2}{return (int)Tokens.STRING_LITERAL_LONG2;}
227  \*                     {return (int)'*';}
228  \{                     {return (int)'{';}
229  \}                     {return (int)'}';}
230  \(                     {return (int)'(';}
231  \)                     {return (int)')';}
232  ,                      {return (int)',';}
233  ;                      {return (int)';';}
234  \.                     {return (int)'.';}
235  a                      {return (int)'a';}
236  \[                     {return (int)'[';}
237  \]                     {return (int)']';}
238  \+                     {return (int)'+';}
239  \-                     {return (int)'-';}
240  \!                     {return (int)'!';}
241  \/                     {return (int)'/';}
242  {ws}                   {/*IGNORE*/}
243  {nil}                  {return (int)Tokens.NIL;}
244  {anon}                 {return (int)Tokens.ANON;}
245  {comment}              {/*IGNORE*/}
246
247  /*
248   * Match for any single character not matching any other token definition.
249   * (Makes parser aware of undefined tokens.)
250   */
251  .                      {return (int)yytext[0];}
252
253  /* Special-case match for end of file. */
254  <<EOF>>                {return (int)Tokens.EOF;}
255
256  %{
257    // Called each time a token has been returned.
258    LoadYylval();
259  %}
260
261  %%
```

# Appendix C

# Parser Specification

```
1    %using SharQL.Ast.Nodes
2    %using SharQL.Ast
3    %using SharQL.Utils
4    %using System.Globalization
5    %using System.Diagnostics
6
7    /* Defines the generated parser namespace. */
8    %namespace SharQL
9
10   /* Makes the generated parser class a partial class. */
11   %partial
12
13   /* Defines the generated parser class name. */
14   %parsertype ParsingEngine
15
16   /* Defines the token location object type name. */
17   %YYLTYPE CustomLexLocation
18
19   %{
20
21   // Provides access to the generated Scanner class instance.
22   private Scanner sc { get { return scanner as Scanner; } }
23
24   %}
25
26   /*
27    * SPARQL Query Language for RDF
28    * W3C Recommendation 15 January 2008
29    *
30    * http://www.w3.org/TR/2008/REC-rdf-sparql-query-20080115/
31    *
32    */
33
34   %start Query
35
36   %token IRI_REF PNAME_NS PNAME_LN BLANK_NODE_LABEL VAR1 VAR2 LANGTAG
37   %token INTEGER DECIMAL DOUBLE INTEGER_POSITIVE DECIMAL_POSITIVE
38   %token DOUBLE_POSITIVE INTEGER_NEGATIVE DECIMAL_NEGATIVE
39   %token DOUBLE_NEGATIVE
40   %token EXPONENT STRING_LITERAL1 STRING_LITERAL2 STRING_LITERAL_LONG1
41   %token STRING_LITERAL_LONG2 ECHAR NIL ANON
42   %token BASE PREFIX SELECT DISTINCT REDUCED CONSTRUCT DESCRIBE ASK FROM
43   %token NAMED WHERE ORDERBY ASC DESC LIMIT OFFSET OPTIONAL
```

```
44  %token GRAPH FILTER
45  %token UNION OR AND NOT LESSOREQUAL GREATEROREQUAL STR LANG LANGMATCHES
46  %token DATATYPE BOUND SAMETERM ISIRI ISURI ISBLANK
47  %token ISLITERAL REGEX TRUE
48  %token FALSE HATS
49
50  %union { public NodeBase Value; }
51
52  %%
53
54  /* Bracketed numbers refer to SPARQL grammar rule numbers.  */
55
56  /* [1] */
57  Query
58  : Prologue SelectQuery
59  {
60    this.AstRoot = new QueryNode($1, $2);
61  }
62  | Prologue ConstructQuery
63  {
64    this.AstRoot = new QueryNode($1, $2);
65  }
66  | Prologue DescribeQuery
67  {
68    this.AstRoot = new QueryNode($1, $2);
69  }
70  | Prologue AskQuery
71  {
72    this.AstRoot = new QueryNode($1, $2);
73  }
74  ;
75
76  /* [2] */
77  Prologue
78  /* REDUCE/REDUCE
79  : BaseDecl
80  {
81    $$ = new PrologueNode($1, null);
82  } */
83  : BaseDecl PrefixDeclList
84  {
85    $$ = new PrologueNode($1, $2);
86  }
87  | PrefixDeclList
88  {
89    $$ = new PrologueNode(null, $1);
90  }
91  /*
92  | REDUCE/REDUCE
93  {
94    $$ = null;
95    Debug.WriteLine("No Prologue");
96  }*/
97  ;
98
99  PrefixDeclList
100 : PrefixDecl PrefixDeclList
101 {
102   $$ = new PrefixDeclListNode($1, $2);
103 }
104 | /* empty */
105 {
```

```
106      $$ = null;
107      Debug.WriteLine("End of PrefixDeclList");
108    }
109    ;
110
111    /* [3] */
112    BaseDecl
113    : BASE IriRef
114    {
115      $$ = new PrefixDeclNode($2);
116    }
117    ;
118
119    /* [4] */
120    PrefixDecl
121    : PREFIX Pname IriRef
122    {
123      PrefixDeclNode pn = (PrefixDeclNode)$2;
124      $$ = new PrefixDeclNode(pn.Namespace, $3);
125    }
126    ;
127
128    Pname
129    : PNAME_NS
130    {
131      $$ = new PrefixDeclNode(sc.yytext, null);
132    }
133    ;
134
135    IriRef
136    : IRI_REF
137    {
138      $$ = new IriRefNode() { Value = sc.yytext };
139    }
140    ;
141
142    /* [5] */
143    SelectQuery
144    : SELECT DistinctReducedModifier StarVarList DatasetClauseList
          WhereClause SolutionModifier
145    {
146      $$ = new SelectQueryNode($2, $3, $4, $5, $6);
147    }
148    ;
149
150    DistinctReducedModifier
151    : DISTINCT
152    {
153      $$ = new DistinctReducedModifierNode()
154      {
155        Modifier = DistinctReducedModifierNode.Modifiers.Distinct
156      };
157    }
158    | REDUCED
159    {
160      $$ = new DistinctReducedModifierNode()
161      {
162        Modifier = DistinctReducedModifierNode.Modifiers.Reduced
163      };
164    }
165    | /* empty */
166    {
```

75

```
167    $$ = null;
168    Debug.WriteLine("No DistinctReducedModified");
169  }
170  ;
171
172  StarVarList
173  : '*'
174  {
175    $$ = new StarVarListNode();
176  }
177  | Var VarList
178  {
179    $$ = new StarVarListNode($1, $2);
180  }
181  ;
182
183  VarList
184  : Var VarList
185  {
186    $$ = new VarListNode($1, $2);
187  }
188  | /* empty */
189  {
190    $$ = null;
191    Debug.WriteLine("Empty VarList");
192  }
193  ;
194
195  DatasetClauseList
196  : DatasetClause DatasetClauseList
197  {
198    $$ = new DatasetClauseListNode($1, $2);
199  }
200  | /* empty */
201  {
202    $$ = null;
203    Debug.WriteLine("Empty DataSetClauseList");
204  }
205  ;
206
207  /* [6] */
208  ConstructQuery
209  : CONSTRUCT ConstructTemplate DatasetClauseList WhereClause
          SolutionModifier
210  {
211    $$ = new ConstructQueryNode($2, $3, $4, $5);
212  }
213  ;
214
215  /* [7] */
216  DescribeQuery
217  : DESCRIBE StarVarOrIRIrefList DatasetClauseList ConditionalWhereClause
          SolutionModifier
218  {
219    $$ = new DescribeQueryNode($2, $3, $4, $5);
220  }
221  ;
222
223  StarVarOrIRIrefList
224  : '*'
225  {
226    $$ = new StarVarOrIRIrefListNode();
```

```
227   }
228   | VarOrIRIref VarOrIRIrefList
229   {
230     $$ = new StarVarOrIRIrefListNode($1, $2);
231   }
232   ;
233
234   VarOrIRIrefList
235   : VarOrIRIref VarOrIRIrefList
236   {
237     $$ = new VarOrIRIrefListNode($1, $2);
238   }
239   | /* empty */
240   {
241     $$ = null;
242     Debug.WriteLine("Empty VarOrIRIrefList");
243   }
244   ;
245
246   ConditionalWhereClause
247   : WhereClause
248   {
249     $$ = new ConditionalWhereClauseNode($1);
250   }
251   | /* empty */
252   {
253     $$ = null;
254     Debug.WriteLine("Empty ConditionalWhereClause");
255   }
256   ;
257
258   /* [8] */
259   AskQuery
260   : ASK DatasetClauseList WhereClause
261   {
262     $$ = new AskQueryNode($2, $3);
263   }
264   ;
265
266   /* [9] */
267   DatasetClause
268   : FROM DefaultOrNamedGraphClause
269   {
270     $$ = new DatasetClauseNode($2);
271   }
272   ;
273
274   DefaultOrNamedGraphClause
275   : DefaultGraphClause
276   {
277     $$ = new DefaultOrNamedGraphClauseNode($1)
278     {
279       Type = DefaultOrNamedGraphClauseNode.Types.Default
280     };
281   }
282   | NamedGraphClause
283   {
284     $$ = new DefaultOrNamedGraphClauseNode($1)
285     {
286       Type = DefaultOrNamedGraphClauseNode.Types.Named
287     };
288   }
```

```
289    ;
290
291    /* [10] */
292    DefaultGraphClause
293    : SourceSelector
294    {
295      $$ = new DefaultGraphClauseNode($1);
296    }
297    ;
298
299    /* [11] */
300    NamedGraphClause
301    : NAMED SourceSelector
302    {
303      $$ = new NamedGraphClauseNode($1);
304    }
305    ;
306
307    /* [12] */
308    SourceSelector
309    : IRIref
310    {
311      $$ = new SourceSelectorNode($1);
312    }
313    ;
314
315    /* [13] */
316    WhereClause
317    : WHERE GroupGraphPattern
318    {
319      $$ = new WhereClauseNode($2);
320    }
321    | GroupGraphPattern
322    {
323      $$ = new WhereClauseNode($1);
324    }
325    ;
326
327    /* [14] */
328    SolutionModifier
329    : OrderClause
330    {
331      $$ = new SolutionModifierNode($1, null);
332    }
333    | LimitOffsetClauses
334    {
335      $$ = new SolutionModifierNode(null, $1);
336    }
337    | OrderClause LimitOffsetClauses
338    {
339      $$ = new SolutionModifierNode($1, $2);
340    }
341    | /* empty */
342    {
343      $$ = null;
344      Debug.WriteLine("Empty SolutionModifier");
345    }
346    ;
347
348    /* [15] */
349    LimitOffsetClauses
350    : LimitClause
```

```
351  {
352    $$ = new LimitOffsetClausesNode($1, null);
353  }
354  | LimitClause OffsetClause
355  {
356    $$ = new LimitOffsetClausesNode($1, $2);
357  }
358  | OffsetClause
359  {
360    $$ = new LimitOffsetClausesNode(null, $1);
361  }
362  | OffsetClause LimitClause
363  {
364    $$ = new LimitOffsetClausesNode($2, $1);
365  }
366  ;
367
368  /* [16] */
369  OrderClause
370  : ORDERBY OrderCondition OrderConditionList
371  {
372    $$ = new OrderClauseNode($2, $3);
373  }
374  ;
375
376  OrderConditionList
377  : OrderCondition OrderConditionList
378  {
379    $$ = new OrderConditionListNode($1, $2);
380  }
381  | /* empty */
382  {
383    $$ = null;
384    Debug.WriteLine("Empty OrderConditionList");
385  }
386  ;
387
388  /* [17] */
389  OrderCondition
390  : AscOrDesc BrackettedExpression
391  {
392    $$ = new OrderConditionNode($1, $2, null, null);
393  }
394  | Constraint
395  {
396    $$ = new OrderConditionNode(null, null, $1, null);
397  }
398  | Var
399  {
400    $$ = new OrderConditionNode(null, null, null, $1);
401  }
402  ;
403
404  AscOrDesc
405  : ASC
406  {
407    $$ = new AscOrDescNode()
408    {
409      Order = AscOrDescNode.Orders.Ascending
410    };
411  }
412  | DESC
```

```
413   {
414     $$ = new AscOrDescNode()
415     {
416       Order = AscOrDescNode.Orders.Descending
417     };
418   }
419   ;
420
421   /* [18] */
422   LimitClause
423   : LIMIT INTEGER
424   {
425     $$ = new LimitClauseNode()
426     {
427       Limit = Convert.ToInt32(sc.yytext)
428     };
429   }
430   ;
431
432   /* [19] */
433   OffsetClause
434   : OFFSET INTEGER
435   {
436     $$ = new OffsetClauseNode()
437     {
438       Offset = Convert.ToInt32(sc.yytext)
439     };
440   }
441   ;
442
443   /* [20] */
444   GroupGraphPattern
445   : '{' ConditionalTriplesBlock GroupGraphPatternBlockList '}'
446   {
447     $$ = new GroupGraphPatternNode($2, $3);
448   }
449   ;
450
451   ConditionalTriplesBlock
452   : TriplesBlock
453   {
454     $$ = $1;
455   }
456   | /* empty */
457   {
458     $$ = null;
459     Debug.WriteLine("No TriplesBlock");
460   }
461   ;
462
463   GroupGraphPatternBlockList
464   : GroupGraphPatternBlock GroupGraphPatternBlockList
465   {
466     $$ = new GroupGraphPatternBlockListNode($1, $2);
467   }
468   | /* empty */
469   {
470     $$ = null;
471     Debug.WriteLine("End of GroupGraphPatternBlockList");
472   }
473   ;
474
```

```
475  GroupGraphPatternBlock
476  : GraphPatternNotTriples ConditionalDot ConditionalTriplesBlock
477  {
478    $$ = new GroupGraphPatternBlockNode($1, null, $3);
479  }
480  | Filter ConditionalDot ConditionalTriplesBlock
481  {
482    $$ = new GroupGraphPatternBlockNode(null, $1, $3);
483  }
484  ;
485
486  ConditionalDot
487  : '.'
488  {
489    Debug.WriteLine("Dot");
490  }
491  | /* empty */
492  {
493    Debug.WriteLine("No Dot");
494  }
495  ;
496
497  /* [21] */
498  TriplesBlock
499  : TriplesSameSubject ConditionalDottedTriplesBlock
500  {
501    $$ = new TriplesBlockNode($1, $2);
502  }
503  ;
504
505  ConditionalDottedTriplesBlock
506  : '.'
507  {
508    $$ = null;
509  }
510  | '.' TriplesBlock
511  {
512    $$ = $2;
513  }
514  | /* empty */
515  {
516    $$ = null;
517    Debug.WriteLine("No ConditionalDottedTriplesBlock");
518  }
519  ;
520
521  /* [22] */
522  GraphPatternNotTriples
523  : OptionalGraphPattern
524  {
525    $$ = $1;
526  }
527  | GroupOrUnionGraphPattern
528  {
529    $$ = $1;
530  }
531  | GraphGraphPattern
532  {
533    $$ = $1;
534  }
535  ;
536
```

81

```
537   /* [23] */
538   OptionalGraphPattern
539   : OPTIONAL GroupGraphPattern
540   {
541     $$ = new OptionalGraphPatternNode($2);
542   }
543   ;
544
545   /* [24] */
546   GraphGraphPattern
547   : GRAPH VarOrIRIref GroupGraphPattern
548   {
549     $$ = new GraphGraphPatternNode($1, $2);
550   }
551   ;
552
553   /* [25] */
554   GroupOrUnionGraphPattern
555   : GroupGraphPattern GroupGraphPatternUnionList
556   {
557     $$ = new GroupOrUnionGraphPatternNode($1, $2);
558   }
559   ;
560
561   GroupGraphPatternUnionList
562   : UNION GroupGraphPattern GroupGraphPatternUnionList
563   {
564     $$ = new GroupGraphPatternUnionListNode($2);
565   }
566   | /* empty */
567   {
568     $$ = null;
569     Debug.WriteLine("End of GroupGraphPatternUnionList");
570   }
571   ;
572
573   /* [26] */
574   Filter
575   : FILTER Constraint
576   {
577     $$ = new FilterNode($2);
578   }
579   ;
580
581   /* [27] */
582   Constraint
583   : BrackettedExpression
584   {
585     $$ = $1;
586   }
587   | BuiltInCall
588   {
589     $$ = $1;
590   }
591   | FunctionCall
592   {
593     $$ = $1;
594   }
595   ;
596
597   /* [28] */
598   FunctionCall
```

```
599    : IRIref ArgList
600    {
601      $$ = new FunctionCallNode($1, $2);
602    }
603    ;
604
605    /* [29] */
606    ArgList
607    : NIL
608    {
609      $$ = new ArgListNode();
610    }
611    | '(' Expression ExpressionList ')'
612    {
613      $$ = new ArgListNode($2, $3);
614    }
615    ;
616
617    ExpressionList
618    : ',' Expression ExpressionList
619    {
620      $$ = new ExpressionListNode($2, $3);
621    }
622    | /* empty */
623    {
624      $$ = null;
625      Debug.WriteLine("End of ExpressionList");
626    }
627    ;
628
629    /* [30] */
630    ConstructTemplate
631    : '{' ConditionalConstructTriples '}'
632    {
633      $$ = $2;
634    }
635    ;
636
637    ConditionalConstructTriples
638    : ConstructTriples
639    {
640      $$ = $1;
641    }
642    | /* empty */
643    {
644      $$ = null;
645      Debug.WriteLine("No ConstructTriples");
646    }
647    ;
648
649    /* [31]  */
650    ConstructTriples
651    : TriplesSameSubject ConditionalConstructDottedTriples
652    {
653      $$ = new ConstructTriplesNode($1, $2);
654    }
655    ;
656
657    ConditionalConstructDottedTriples
658    : '.'
659    {
660      $$ = new ConditionalConstructDottedTriplesNode(null);
```

```
661  }
662  | '.' ConstructTriples
663  {
664    $$ = new ConditionalConstructDottedTriplesNode($2);
665  }
666  | /* empty */
667  {
668    $$ = null;
669    Debug.WriteLine("No ConstructTriples");
670  }
671  ;
672
673  /* [32] */
674  TriplesSameSubject
675  : VarOrTerm PropertyListNotEmpty
676  {
677    $$ = new TriplesSameSubjectNode($1, null, $2);
678  }
679  | TriplesNode PropertyList
680  {
681    $$ = new TriplesSameSubjectNode(null, $1, $2);
682  }
683  ;
684
685  /* [33] */
686  PropertyListNotEmpty
687  /* REDUCE/REDUCE
688  : Verb ObjectList
689  {
690    $$ = new PropertyListNode($1, $2, null);
691  }*/
692  : Verb ObjectList PropertyListTail
693  {
694    $$ = new PropertyListNode($1, $2, $3);
695  }
696  ;
697
698  PropertyListTail
699  : ';' PropertyListTail
700  {
701    $$ = $2;
702  }
703  | ';' Verb ObjectList PropertyListTail
704  {
705    $$ = new PropertyListNode($2, $3, $4);
706  }
707  | /* empty */
708  {
709    $$ = null;
710    Debug.WriteLine("End of PropertyList");
711  }
712  ;
713
714  /* [34] */
715  PropertyList
716  : PropertyListNotEmpty
717  {
718    $$ = $1;
719  }
720  | /* empty */
721  {
722    $$ = null;
```

```
723      Debug.WriteLine("End of PropertyList");
724    }
725    ;
726
727    /* [35] */
728    ObjectList
729    /* REDUCE/REDUCE
730    : Object
731    {
732      $$ = new ObjectListNode($1, null);
733    } */
734    : Object ObjectListTail
735    {
736      $$ = new ObjectListNode($1, $2);
737    }
738    ;
739
740    ObjectListTail
741    : ',' Object ObjectListTail
742    {
743      $$ = new ObjectListNode($2, $3);
744    }
745    | /* empty */
746    {
747      $$ = null;
748      Debug.WriteLine("End of ObjectListTail");
749    }
750    ;
751
752    /* [36] */
753    Object
754    : GraphNode
755    {
756      $$ = $1;
757    }
758    ;
759
760    /* [37] */
761    Verb
762    : VarOrIRIref
763    {
764      $$ = $1;
765    }
766    | 'a'
767    {
768      $$ = new IriRefNode()
769      {
770        Value = "http://www.w3.org/1999/02/22-rdf-syntax-ns#type"
771      };
772    }
773    ;
774
775    /* [38] */
776    TriplesNode
777    : Collection
778    {
779      $$ = $1;
780    }
781    | BlankNodePropertyList
782    {
783      $$ = $1;
784    }
```

85

```
785   ;
786
787   /* [39] */
788   BlankNodePropertyList
789   : '[' PropertyListNotEmpty ']'
790   {
791     $$ = new BlankNodePropertyListNode($2);
792   }
793   ;
794
795   /* [40] */
796   Collection
797   : '(' GraphNode GraphNodeList ')'
798   {
799     $$ = new CollectionNode(new GraphNodeListNode($2, $3));
800   }
801   ;
802
803   GraphNodeList
804   : GraphNode GraphNodeList
805   {
806     $$ = new GraphNodeListNode($1, $2);
807   }
808   | /* empty */
809   {
810     $$ = null;
811     Debug.WriteLine("End of GraphNodeList");
812   }
813   ;
814
815   /* [41] */
816   GraphNode
817   : VarOrTerm
818   {
819     $$ = $1;
820   }
821   | TriplesNode
822   {
823     $$ = $1;
824   }
825   ;
826
827   /* [42] */
828   VarOrTerm
829   : Var
830   {
831     $$ = $1;
832   }
833   | GraphTerm
834   {
835     $$ = $1;
836   }
837   ;
838
839   /* [43] */
840   VarOrIRIref
841   : Var
842   {
843     $$ = $1;
844   }
845   | IRIref
846   {
```

```
847       $$ = $1;
848    }
849    ;
850
851    /* [44] */
852    Var
853    : VAR1
854    {
855       $$ = new VarNode(sc.yytext.Substring(1));
856    }
857    | VAR2
858    {
859       $$ = new VarNode(sc.yytext.Substring(1));
860    }
861    ;
862
863    /* [45] */
864    GraphTerm
865    : IRIref
866    {
867       $$ = $1;
868    }
869    | RDFLiteral
870    {
871       $$ = $1;
872    }
873    | NumericLiteral
874    {
875       $$ = $1;
876    }
877    | BooleanLiteral
878    {
879       $$ = $1;
880    }
881    | BlankNode
882    {
883       $$ = $1;
884    }
885    | NIL
886    {
887       $$ = null;
888    }
889    ;
890
891    /* [46] */
892    Expression
893    : ConditionalOrExpression
894    {
895       $$ = $1;
896    }
897    ;
898
899    /* [47] */
900    ConditionalOrExpression
901    : ConditionalAndExpression ConditionalOrExpressionList
902    {
903       $$ = new ConditionalOrExpressionNode($1, $2);
904    }
905    ;
906
907    ConditionalOrExpressionList
908    : OR ConditionalAndExpression ConditionalOrExpressionList
```

```
909  {
910    $$ = new ConditionalOrExpressionListNode($2, $3);
911  }
912  | /* empty */
913  {
914    $$ = null;
915    Debug.WriteLine("End of ConditionalOrExpressionList");
916  }
917  ;
918
919  /* [48] */
920  ConditionalAndExpression
921  : ValueLogical ConditionalAndExpressionList
922  {
923    $$ = new ConditionalAndExpressionNode($1, $2);
924  }
925  ;
926
927  ConditionalAndExpressionList
928  : AND ValueLogical ConditionalAndExpressionList
929  {
930    $$ = new ConditionalAndExpressionListNode($2, $3);
931  }
932  | /* empty */
933  {
934    $$ = null;
935    Debug.WriteLine("End of ConditionalAndExpressionList");
936  }
937  ;
938
939  /* [49] */
940  ValueLogical
941  : RelationalExpression
942  {
943    $$ = $1;
944  }
945  ;
946
947  /* [50] */
948  RelationalExpression
949  : NumericExpression
950  {
951    $$ = $1;
952  }
953  | NumericExpression '=' NumericExpression
954  {
955    $$ = new RelationalExpressionNode(Constants.Operators.Equal, $1, $3);
956  }
957  | NumericExpression NOT NumericExpression
958  {
959    $$ = new RelationalExpressionNode(Constants.Operators.Not, $1, $3);
960  }
961  | NumericExpression '<' NumericExpression
962  {
963    $$ = new RelationalExpressionNode(Constants.Operators.Less, $1, $3);
964  }
965  | NumericExpression '>' NumericExpression
966  {
967    $$ = new RelationalExpressionNode(Constants.Operators.Greater, $1, $3);
968  }
969  | NumericExpression LESSOREQUAL NumericExpression
970  {
```

```
971    $$ = new RelationalExpressionNode(Constants.Operators.LessOrEqual, $1,
          $3);
972  }
973  | NumericExpression GREATEROREQUAL NumericExpression
974  {
975    $$ = new RelationalExpressionNode(Constants.Operators.GreaterOrEqual,
          $1, $3);
976  }
977  ;
978
979  /* [51] */
980  NumericExpression
981  : AdditiveExpression
982  {
983    $$ = $1;
984  }
985  ;
986
987  /* [52] */
988  AdditiveExpression
989  : MultiplicativeExpression AdditiveExpressionList
990  {
991    $$ = new AdditiveExpressionNode($1, $2);
992  }
993  ;
994
995  AdditiveExpressionList
996  : '+' MultiplicativeExpression AdditiveExpressionList
997  {
998    $$ = new AdditiveExpressionListNode(Constants.Operators.Add, $1, $2);
999  }
1000 | '-' MultiplicativeExpression AdditiveExpressionList
1001 {
1002   $$ = new AdditiveExpressionListNode(Constants.Operators.Subtract, $1,
          $2);
1003 }
1004 | NumericLiteralPositive AdditiveExpressionList
1005 {
1006   $$ = new AdditiveExpressionListNode(Constants.Operators.None, $1, $2);
1007 }
1008 | NumericLiteralNegative AdditiveExpressionList
1009 {
1010   $$ = new AdditiveExpressionListNode(Constants.Operators.None, $1, $2);
1011 }
1012 | /* empty */
1013 {
1014   $$ = null;
1015   Debug.WriteLine("End of AdditiveExpressionList");
1016 }
1017 ;
1018
1019 /* [53] */
1020 MultiplicativeExpression
1021 : UnaryExpression UnaryExpressionList
1022 {
1023   $$ = new MultiplicativeExpressionNode($1, $2);
1024 }
1025 ;
1026
1027 UnaryExpressionList
1028 : '*' UnaryExpression UnaryExpressionList
1029 {
```

```
1030      $$ = new UnaryExpressionListNode(Constants.Operators.Multiply, $2, $3);
1031    }
1032    | '/' UnaryExpression UnaryExpressionList
1033    {
1034      $$ = new UnaryExpressionListNode(Constants.Operators.Divide, $2, $3);
1035    }
1036    | /* empty */
1037    {
1038      $$ = null;
1039      Debug.WriteLine("End of UnaryExpressionList");
1040    }
1041    ;
1042
1043    /* [54] */
1044    UnaryExpression
1045    : '!' PrimaryExpression
1046    {
1047      $$ = new UnaryExpressionNode(Constants.Operators.Not, $2);
1048    }
1049    | '+' PrimaryExpression
1050    {
1051      $$ = new UnaryExpressionNode(Constants.Operators.Add, $2);
1052    }
1053    | '-' PrimaryExpression
1054    {
1055      $$ = new UnaryExpressionNode(Constants.Operators.Subtract, $2);
1056    }
1057    | PrimaryExpression
1058    {
1059      $$ = $1;
1060    }
1061    ;
1062
1063    /* [55] */
1064    PrimaryExpression
1065    : BrackettedExpression
1066    {
1067      $$ = $1;
1068    }
1069    | BuiltInCall
1070    {
1071      $$ = $1;
1072    }
1073    | IRIrefOrFunction
1074    {
1075      $$ = $1;
1076    }
1077    | RDFLiteral
1078    {
1079      $$ = $1;
1080    }
1081    | NumericLiteral
1082    {
1083      $$ = $1;
1084    }
1085    | BooleanLiteral
1086    {
1087      $$ = $1;
1088    }
1089    | Var
1090    {
1091      $$ = $1;
```

```
1092  }
1093  ;
1094
1095  /* [56] */
1096  BrackettedExpression
1097  : '(' Expression ')'
1098  {
1099    $$ = $2;
1100  }
1101  ;
1102
1103  /* [57] */
1104  BuiltInCall
1105  : STR '(' Expression ')'
1106  {
1107    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.Str, $3);
1108  }
1109  | LANG '(' Expression ')'
1110  {
1111    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.Lang, $3);
1112  }
1113  | LANGMATCHES '(' Expression ',' Expression ')'
1114  {
1115    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.LangMatches, $3, $5);
1116  }
1117  | DATATYPE '(' Expression ')'
1118  {
1119    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.Datatype, $3);
1120  }
1121  | BOUND '(' Var ')'
1122  {
1123    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.Bound, $3);
1124  }
1125  | SAMETERM '(' Expression ',' Expression ')'
1126  {
1127    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.SameTerm, $3, $5);
1128  }
1129  | ISIRI '(' Expression ')'
1130  {
1131    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.IsIri, $3);
1132  }
1133  | ISURI '(' Expression ')'
1134  {
1135    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.IsUri, $3);
1136  }
1137  | ISBLANK '(' Expression ')'
1138  {
1139    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.IsBlank, $3);
1140  }
1141  | ISLITERAL '(' Expression ')'
1142  {
1143    $$ = new BuiltInCallNode(BuiltInCallNode.CallType.IsLiteral, $3);
1144  }
1145  | RegexExpression
1146  {
1147    $$ = $1;
1148  }
1149  ;
1150
1151  /* [58] */
1152  RegexExpression
1153  : REGEX '(' Expression ',' Expression ')'
```

```
1154  {
1155    $$ = new RegexExpressionNode($3, $5);
1156  }
1157  | REGEX '(' Expression ',' Expression ',' Expression ')'
1158  {
1159    $$ = new RegexExpressionNode($3, $5, $7);
1160  }
1161  ;
1162
1163  /* [59] */
1164  IRIrefOrFunction
1165  : IRIref
1166  {
1167    $$ = $1;
1168  }
1169  | IRIref ArgList
1170  {
1171    $$ = new FunctionCallNode($1, $2);
1172  }
1173  ;
1174
1175  /* [60] */
1176  RDFLiteral
1177  : String
1178  {
1179    $$ = new RdfLiteralNode($1, null);
1180  }
1181  | String LANGTAG
1182  {
1183    $$ = new RdfLiteralNode($1, null) { LangTag = sc.yytext };
1184  }
1185  | String HATS IRIref
1186  {
1187    $$ = new RdfLiteralNode($1, $3);
1188  }
1189  ;
1190
1191  /* [61] */
1192  NumericLiteral
1193  : NumericLiteralUnsigned
1194  {
1195    $$ = $1;
1196  }
1197  | NumericLiteralPositive
1198  {
1199    $$ = $1;
1200  }
1201  | NumericLiteralNegative
1202  {
1203    $$ = $1;
1204  }
1205  ;
1206
1207  /* [62] */
1208  NumericLiteralUnsigned
1209  : INTEGER
1210  {
1211    $$ = new NumericLiteralNode() { Type =
1212        NumericLiteralNode.LiteralType.Int,
1213      IntValue = Convert.ToInt32(sc.yytext) };
1213  }
1214  | DECIMAL
```

92

```
1215  {
1216    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Decimal,
1217      DecimalValue = Convert.ToDecimal(sc.yytext,
             CultureInfo.InvariantCulture) };
1218  }
1219  | DOUBLE
1220  {
1221    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Double,
1222      DoubleValue = Convert.ToDouble(sc.yytext,
             CultureInfo.InvariantCulture) };
1223  }
1224  ;
1225
1226  /* [63] */
1227  NumericLiteralPositive
1228  : INTEGER_POSITIVE
1229  {
1230    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Int,
1231      IntValue = Convert.ToInt32(sc.yytext) };
1232  }
1233  | DECIMAL_POSITIVE
1234  {
1235    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Decimal,
1236      DecimalValue = Convert.ToDecimal(sc.yytext,
             CultureInfo.InvariantCulture) };
1237  }
1238  | DOUBLE_POSITIVE
1239  {
1240    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Double,
1241      DoubleValue = Convert.ToDouble(sc.yytext,
             CultureInfo.InvariantCulture) };
1242  }
1243  ;
1244
1245  /* [64] */
1246  NumericLiteralNegative
1247  : INTEGER_NEGATIVE
1248  {
1249    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Int,
1250      IntValue = Convert.ToInt32(sc.yytext) };
1251  }
1252  | DECIMAL_NEGATIVE
1253  {
1254    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Decimal,
1255      DecimalValue = Convert.ToDecimal(sc.yytext,
             CultureInfo.InvariantCulture) };
1256  }
1257  | DOUBLE_NEGATIVE
1258  {
1259    $$ = new NumericLiteralNode() { Type =
             NumericLiteralNode.LiteralType.Double,
1260      DoubleValue = Convert.ToDouble(sc.yytext,
             CultureInfo.InvariantCulture) };
1261  }
1262  ;
```

```
1263
1264   /* [65] */
1265   BooleanLiteral
1266   : TRUE
1267   {
1268     $$ = new BooleanLiteralNode() { Value = true };
1269   }
1270   | FALSE
1271   {
1272     $$ = new BooleanLiteralNode() { Value = false };
1273   }
1274   ;
1275
1276   /* [66] */
1277   String
1278   : STRING_LITERAL1
1279   {
1280     string str = sc.yytext;
1281     $$ = new StringLiteralNode() { Value = str.Substring(1, str.Length - 2)
          };
1282   }
1283   | STRING_LITERAL2
1284   {
1285     string str = sc.yytext;
1286     $$ = new StringLiteralNode() { Value = str.Substring(1, str.Length - 2)
          };
1287   }
1288   | STRING_LITERAL_LONG1
1289   {
1290     string str = sc.yytext;
1291     $$ = new StringLiteralNode() { Value = str.Substring(3, str.Length - 6)
          };
1292   }
1293   | STRING_LITERAL_LONG2
1294   {
1295     string str = sc.yytext;
1296     $$ = new StringLiteralNode() { Value = str.Substring(3, str.Length - 6)
          };
1297   }
1298   ;
1299
1300   /* [67] */
1301   IRIref
1302   : IRI_REF
1303   {
1304     $$ = new IriRefNode() { Value = sc.yytext };
1305   }
1306   | PrefixedName
1307   {
1308     $$ = $1;
1309   }
1310   ;
1311
1312   /* [68] */
1313   PrefixedName
1314   : PNAME_LN
1315   {
1316     string str = sc.yytext;
1317     $$ = new PrefixedNameNode() { Pname_Ns = str.Split(':')[0], Pn_Local =
          str.Substring(str.IndexOf(':') + 1) };
1318   }
1319   | PNAME_NS
```

```
1320   {
1321     $$ = new PrefixedNameNode() { Pname_Ns = sc.yytext.Split(':')[0] };
1322   }
1323   ;
1324
1325   /* [69] */
1326   BlankNode
1327   : BLANK_NODE_LABEL
1328   {
1329     $$ = new BlankNode(){ Value = sc.yytext };
1330   }
1331   | ANON
1332   {
1333     $$ = new BlankNode();
1334   }
1335   ;
1336
1337   %%
```

# Appendix D

# *NodeBase* Class

```csharp
1  using System.Collections.Generic;
2  using SharQL.Ast.Visitor;
3  using System.Diagnostics;
4
5  namespace SharQL.Ast
6  {
7    /// <summary>
8    /// The base implementation of the <see cref="INode"/> interface.
9    /// </summary>
10   public abstract class NodeBase : INode
11   {
12     #region Protected Fields
13
14     /// <summary>The parent of the node.</summary>
15     protected INode parent;
16     /// <summary>The children of the node.</summary>
17     protected IList<INode> children;
18
19     #endregion
20
21     /// <summary>
22     /// Gets a string uniquely identifying the type of the node.
23     /// </summary>
24     /// <remarks>Obtains the type name by reflection. This can
25     /// be slow. For better performance, override and return a
26     /// constant.</remarks>
27     public virtual string NodeType
28     {
29       get { return this.GetType().FullName; }
30     }
31
32     /// <summary>
33     /// Default constructor.
34     /// </summary>
35     public NodeBase()
36     {
37       children = new NodeCollection();
38     }
39
40     /// <summary>
41     /// Constructs the object and adds a set of children.
42     /// </summary>
43     /// <param name="children">The children to add</param>
```

```csharp
44        public NodeBase(params NodeBase[] children)
45          : this()
46        {
47          foreach (NodeBase node in children)
48          {
49            if(node != null) node.Parent = this;
50            this.children.Add(node);
51          }
52        }
53
54        /// <summary>
55        /// Gets or sets the parent of the node. Should be set to
56        /// <c>null</c> if node is root.
57        /// </summary>
58        public virtual INode Parent
59        {
60          get { return parent; }
61          set { parent = value; }
62        }
63
64        /// <summary>
65        /// Gets the list of children of the node.
66        /// </summary>
67        public virtual IList<INode> Children
68        {
69          get { return children; }
70        }
71
72        /// <summary>
73        /// Accepts a visitor. The visitor pattern can be used for
74        /// e.g. optimizing the AST.
75        /// </summary>
76        /// <param name="visitor">The visitor</param>
77        public virtual void Accept(IParserVisitor visitor)
78        {
79          foreach (INode child in children)
80          {
81            if (child != null)
82              child.Accept(visitor);
83          }
84          visitor[NodeType](this);
85        }
86
87        /// <summary>
88        /// Returns a string representation of the node.
89        /// </summary>
90        /// <returns>A string representation of the node</returns>
91        public override string ToString()
92        {
93          return NodeType;
94        }
95
96        /// <summary>
97        /// Returns a prefixed string representation of the node.
98        /// </summary>
99        /// <param name="prefix">The string prefix to which the string
100       /// representation should be appended.</param>
101       /// <returns>A prefixed string representation of the node</returns>
102       public virtual string ToString(string prefix)
103       {
104         return prefix + ToString();
105       }
```

```
106          }
107      }
```

# Appendix E

# The W3C SPARQL Test Suite

The tests included in the W3C SPARQL Test Suite are listed in the following tables.

| syntax-sparql1 | |
|---|---|
| **Category** | **Tests included** |
| **Basic** | `syntax-basic-01...syntax-basic-06` |
| **Blank Nodes** | `syntax-bnodes-01...syntax-bnodes-05` |
| **Expressions** | `syntax-expr-01...syntax-expr-05` |
| **Forms** | `syntax-forms-01, syntax-forms-02` |
| **Limit/Offset** | `syntax-limit-offset-01...` `syntax-limit-offset-04` |
| **Lists** | `syntax-lists-01...syntax-lists-05` |
| **Literals** | `syntax-lit-01...syntax-lit-20` |
| **Ordering** | `syntax-order-01...syntax-order-07` |
| **Patterns** | `syntax-pat-01...syntax-pat-04` |
| **Naming** | `syntax-qname-01...syntax-qname-08` |
| **Structural** | `syntax-struct-01...syntax-struct-14` |
| **Union** | `syntax-union-01, syntax-union-01` |

Table E.1: Tests included in the W3C Test Suite part 1

| syntax-sparql2 | |
|---|---|
| **Category** | **Tests included** |
| **Blank Nodes** | `syntax-bnode-01...syntax-bnode-03` |
| **Dataset** | `syntax-dataset-01...syntax-dataset-04` |
| **Escapes** | `syntax-esc-01...syntax-esc-05` |
| **Forms** | `syntax-form-ask-02,`<br>`syntax-form-construct01...`<br>`syntax-form-construct06,`<br>`syntax-form-describe01,`<br>`syntax-form-describe02,`<br>`syntax-form-select01,`<br>`syntax-form-select02` |
| **Function** | `syntax-function-01...syntax-function-04` |
| **General** | `syntax-general-01...syntax-general-14` |
| **Graphs** | `syntax-graph-01...syntax-graph-05` |
| **Keywords** | `syntax-keywords-01...syntax-keywords-03` |
| **Lists** | `syntax-lists-01...syntax-lists-05` |

Table E.2: Tests included in the W3C Test Suite part 2

| syntax-sparql3 | |
|---|---|
| **Category** | **Tests included** |
| **Positive Tests** | `syn-01...syn-08,` `syn-blabel-cross-filter` |
| **Negative Tests** | `syn-bad-01...syn-bad-31,`<br>`syn-bad-bnode-dot,`<br>`syn-bad-bnodes-missing-pvalues-01,`<br>`syn-bad-bnodes-missing-pvalues-02,`<br>`syn-bad-empty-optional-01,`<br>`syn-bad-empty-optional-02,`<br>`syn-bad-filter-missing-parens,`<br>`syn-bad-lone-list,` `syn-bad-lone-node,`<br>`syn-blabel-cross-graph-bad,`<br>`syn-blabel-cross-optional-bad,`<br>`syn-blabel-cross-union-bad` |

Table E.3: Tests included in the W3C Test Suite part 3

| syntax-sparql4 | |
|---|---|
| **Category** | **Tests included** |
| **Positive Tests** | `syn-09...syn-11,`<br>`syn-leading-digits-in-prefixed-names` |
| **Negative Tests** | `syn-bad-34...syn-bad-38,`<br>`syn-bad-GRAPH-breaks-BGP,`<br>`syn-bad-OPT-breaks-BGP,`<br>`syn-bad-UNION-breaks-BGP` |

Table E.4: Tests included in the W3C Test Suite part 4

| syntax-sparql5 | |
|---|---|
| **Category** | **Tests included** |
| **Reduced** | `syntax-reduced-01, syntax-reduced-01` |

Table E.5: Tests included in the W3C Test Suite part 5

# Appendix F

# Source Code and Report

On the enclosed CD-ROM, the complete source code for the SharQL parser is included, as well as a digital copy of this report.