

**AUTOMATIC GENERATION OF MULTI-LANGUAGE  
OBJECT DOMAIN MODELS THROUGH A SHAPE  
EXPRESSIONS SUBSET**

GUILLERMO FACUNDO COLUNGA

A THESIS SUBMITTED FOR THE DEGREE OF SOFTWARE ENGINEER

DEPARTMENT OF COMPUTER SCIENCE  
UNIVERSITY OF OVIEDO

2020

## Acknowledgments

Write your acknowledgments here.

## Abstract

Surface integration is an important step for automatic 3D reconstruction of real objects. The goal of a surface integration algorithm is to reconstruct a surface from a set of range images registered in a common coordinate system. Based on the surface representation used, existing algorithms can be divided into two categories: volume-based and mesh-based. Volume-based methods have been shown to be robust to scanner noise and small features (regions of high curvature) and can build water tight models of high quality. It is, however, difficult to choose the appropriate voxel size when the input range images have both small features and large registration errors compared to the sampling density of range images. Mesh-based methods are more efficient and need less memory compared to volume-based methods but these methods fail in the presence of small features and are not robust to scanning noise.

This paper presents a robust algorithm for mesh-based surface integration of a set of range images. The algorithm is incremental and operates on a range image and the model reconstructed so far. Our algorithm first, transform the model in the coordinate system of the range image. Then, it finds the regions of model overlapping with the range image. This is done by shooting rays from the scanner, through the vertices in the range image and intersecting them with the model. Finally, the algorithm integrates the overlapping regions by using weighted average of points in the model and the range image. The weights are computed using the scanner uncertainty and helps in reducing the effects of scanning noise. To handle small features robustly the integration of overlapping regions is done by computing the position of vertices in the range image along the scanner's line of sight. Since for every point in a range image there is exactly one depth value, the reconstructed surface in the regions of high curvature will not have self-intersections.

# Contents

<b>List of Figures</b> . . . . .	<b>vii</b>
<b>List of Tables</b> . . . . .	<b>ix</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Motivation . . . . .	1
1.2 Contributions . . . . .	3
1.3 Structure of the Document . . . . .	4
<b>2 Teoretical Background</b>	<b>5</b>
2.1 RDF . . . . .	5
2.2 Validating RDF . . . . .	6
2.2.1 Shape Expressions . . . . .	7
2.2.2 Other Technologies . . . . .	10
2.3 Programming Languages . . . . .	10
2.4 Compilers . . . . .	10
2.4.1 Internal Structure . . . . .	11
2.4.2 Conventional Compilers . . . . .	12
2.4.3 Modern Compilers . . . . .	12
<b>3 Related Work</b>	<b>13</b>
3.1 Simplifications of ShEx . . . . .	13
3.1.1 The <b>S</b> language . . . . .	13
3.1.2 ShExJ Micro Spec . . . . .	13
3.2 ShEx Ecosystem Tools . . . . .	14
3.2.1 Validators . . . . .	14
3.2.2 IDEs . . . . .	15
3.2.3 Others . . . . .	15
<b>I Enhancing Error and Warning Detection and Emission on ShEx</b>	<b>17</b>
<b>4 Analysis of Existing Sintactic and Semantic Analizers</b>	<b>18</b>
4.1 Methodology . . . . .	18
4.2 Sintactic Analyzers . . . . .	19
4.3 Semantic Analyzers . . . . .	20
4.4 Possible Enhancements . . . . .	21

<b>5</b>	<b>Proposed Sintactic and Semantic Analyzer</b>	<b>22</b>
5.1	Structure . . . . .	22
5.1.1	Parser . . . . .	22
5.1.2	Sintactic Analyzer . . . . .	23
5.1.3	Semantic Analyzer . . . . .	23
5.2	Implementation . . . . .	24
5.2.1	Parser . . . . .	24
5.2.2	Sintactic Analizer . . . . .	25
5.2.3	Semantic Analyzer . . . . .	26
5.3	Sintactic and Semantic Error and Warnings Detected . . . . .	27
5.3.1	Not trailing semicolon at last triple constraint . . . . .	27
5.3.2	Prefix not defined . . . . .	27
5.3.3	Shape not defined . . . . .	27
5.3.4	Prefix overridden . . . . .	27
5.3.5	Shape overridden . . . . .	27
5.3.6	Unused prefix definition . . . . .	27
5.3.7	Base set but not used . . . . .	27
<b>II</b>	<b>Translating ShEx Schemas to Object Domain Models</b>	<b>28</b>
<b>6</b>	<b>Object Domain Model Translation Problem</b>	<b>29</b>
6.1	Shape Expressions Expressivity . . . . .	29
6.1.1	Properties Set Defined Shapes . . . . .	30
6.1.2	Property Node Constraints . . . . .	31
6.2	Plain Objects Expressivity . . . . .	31
6.2.1	Plain Objects Structure . . . . .	31
6.2.2	Plain Objects Language Expressivity Dependance . . . . .	32
6.2.3	Plain Objects Expressivity Generalization . . . . .	33
6.3	Shape Expressions and Plain Objects Expressivity Comparison . . . . .	33
<b>7</b>	<b>Proposed Translator</b>	<b>34</b>
7.1	Structure . . . . .	35
7.2	Generated Obejcts . . . . .	35
<b>III</b>	<b>Project Sinthesis</b>	<b>36</b>
<b>8</b>	<b>Evaluation of Results</b>	<b>37</b>
8.1	Methodology . . . . .	37
<b>9</b>	<b>Planning and Budget</b>	<b>38</b>

---

9.1	Planning . . . . .	38
9.2	Budget . . . . .	38
<b>10</b>	<b>Conclusions</b>	<b>39</b>
10.1	Future Work . . . . .	39
<b>IV</b>	<b>Annexes and References</b>	<b>40</b>
<b>Appendix A</b>	<b>ShEx Micro Language</b>	<b>41</b>
A.1	Syntax Specification . . . . .	41
A.2	Lexical Specification . . . . .	41
<b>Appendix B</b>	<b>ShEx-Lite Antlr Grammar</b>	<b>42</b>
B.1	Syntax Specification . . . . .	42
B.2	Lexical Specification . . . . .	42
<b>Appendix C</b>	<b>Project Communications</b>	<b>43</b>
C.1	Open Source Community . . . . .	43
C.2	Scientific Disclosure . . . . .	43
C.3	Dublin Core Meetings . . . . .	43
	<b>References</b>	<b>44</b>

## List of Figures

1.1	The 5 star steps of Linked Data . . . . .	2
2.1	RDF N-Triples Example . . . . .	5
2.2	RDF N-Triples Graph Example . . . . .	6
2.3	RDF Example graph . . . . .	6
2.4	RDF node and its shape . . . . .	6
2.5	Shape Expression Example . . . . .	7
2.6	Shapes, shape expression labels and triple expressions . . . . .	8
2.7	Parts of a triple expression . . . . .	8
2.8	Compiler stages . . . . .	11
3.1	ShEx-Lite integration with Shexer . . . . .	16
4.1	Examples of ShEx micro Compact Syntax code containing syntactic and semantic errors or warnings . . . . .	19
5.1	Syntactic and Semantic Analyzer structure . . . . .	22
5.2	Syntax Tree twenty first nodes produced by the parser . . . . .	23
5.3	Abstract Syntax Tree produced after validation and transformations . . . . .	24
5.4	Checker implementation for missing semicolons warning generation . . . . .	25
5.5	Syntactic warning produced by the proposed syntactic analyzer . . . . .	26
5.6	Common information stored at any AST node . . . . .	26
6.1	Schema modeling a <b>Person</b> in <b>ShExC</b> syntax to the left. And the expected translated code in <b>Java</b> to the right. . . . .	29
6.2	ShEx data model. . . . .	30
6.3	ShEx example schemas. . . . .	30
6.4	Property set defined shape model. . . . .	31
6.5	Java, Python and Rust codings of Person object. . . . .	31
6.6	Java plain object decomposition. . . . .	32
6.7	Rust struct modeling a <b>Person</b> to the left. And the most similar approximation in Java to te right. In the Java approximation the Pet class is an interface that it is inherited by the Cat and Dog classes, that way we allow to store in the variable <b>owningPet</b> values of type Cat and Dog. . . . .	33
7.1	Mental model of ShEx-Lite in the existing ShEx syntaxes context. . . . .	34

---

7.2	Constraints and checks context diagram for ShEx-Lite and ShEx. . . . .	35
7.3	Schema modeling a <b>Person</b> in <b>shexl</b> syntax to the left. And the <b>ShEx-Lite</b> generated code in <b>Java</b> to the right. . . . .	35



## List of Tables

4.1	Detection of the different sintactic errors by the current existing ShEx tools that sintactically analyze the shape expressions. . . . .	20
4.2	Detection of the different semantic errors by the current existing ShEx tools that semantically analyze the shape expressions. . . . .	20

# CHAPTER 1

## Introduction

This chapter covers the motivation, the contributions and the structure of the document. The aim for this chapter is that after reading it the reader already have general idea about the motivation of this dissertation, what it covers and its contributions.

### 1.1 Motivation

Each day more and more devices generate data both automatically and manually, and also each day the development of application in different domains that are backed by databases and expose these data to the web becomes easier. The amount and diversity of data produced clearly exceeds our capacity to consume it.

To describe the data that is so large and complex that traditional data processing applications can't handle the term big-data [12] has emerged. Big data has been described by at least three words starting by V: volume, velocity, variety. Although volume and velocity are the most visible features, variety is a key concept which prevents data integration and generates lots of interoperability problems.

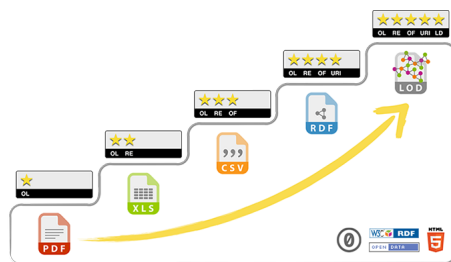
In order to solve this key concept RDF (*Resource Description Framework*) was proposed as a graph-based data model [20] which became part of the Semantic Web [10] vision. Its reliance on the global nature of URIs<sup>1</sup> offered a solution to the data integration problem as RDF datasets produced by different means can seamlessly be integrated with other data.

Also, and related to this is the concept of Linked Data [17] that was proposed as a set of best practices to publish data on the Web. It was introduced by Tim Berners-Lee and was based on four main principles:

- Use URIs as names for things.
- Use HTTP URIs so that people can look up those names.

---

<sup>1</sup>A Uniform Resource Identifier (URI) is a string of characters that unambiguously identifies a particular resource. To guarantee uniformity, all URIs follow a predefined set of syntax rules, but also maintain extensibility through a separately defined hierarchical naming scheme. Ref.[https://en.wikipedia.org/wiki/Uniform\\_Resource\\_Identifier](https://en.wikipedia.org/wiki/Uniform_Resource_Identifier)



**Figure 1.1:** The 5 star steps of Linked Data.

- When someone looks up a URI, provide useful information, using the standards (RDF, SPARQL).
- Include links to other URIs. so that they can discover more things.

This four principles are called the 5 stars Linked Open Data Model, illustrated in [Figure 1.1](#). RDF is mentioned in the third principle as one of the standards that provides useful information. The goal of this principles is that data is not only ready for humans to navigate through but also for other agents, like computers, that may automatically process that data.

All the above motivations helped to make RDF the language for the Web of Data, as described in [\[19\]](#). And the main features that it presents are: *Disambiguation*, *Integration*, *Extensibility*, *Flexibility* and *Open by Default*. With the features also some drawbacks are associated, the most important one and the one we will focus is the **RDF production/consumption dilemma**.

RDF production/consumption dilemma states that it is necessary to find ways that data producers can generate their data so it can be handled by potential consumers. For example, they may want to declare that some nodes have some properties with some specific values. Data consumers need to know that structure to develop applications to consume the data.

Although RDF is a very flexible schema-less language, enterprise and industrial applications may require an extra level of validation before processing for several reasons like security, performance, etc.

To solve that dilemma and as an alternative to expecting the data to have some structure without validation, ShEx (*Shape Expressions*) where proposed as a human-readable and high-level open source language for RDF validation. Initially ShEx was proposed as a human-readable syntax for OSLC Resource Shapes [\[24\]](#) but ShEx grew very fast to embrace more complex user requirements coming from clinical and library use cases.

Another technology, SPIN, was used for RDF validation, principally in TopQuadrant's TopBraid Composer. This technology, influenced from OSLC Resource Shapes as well, evolved into both a private implementation and open source definition of the SHACL (*Shapes*

*Constraint Language*), which was adopted by the W3C Data Shapes Working Group.

From a user point of view the possibilities of ShEx are very large, from the smallest case to just validate a node with one property to a scientific domain case where we need to validate the human genome (*a real use case of ShEx*). A language with such a number of possibilities requires from a strong syntactic and semantic validation and that leads us to our first research question.

**Research Question 1.** *How much the existing syntactic and semantic validation systems for shape expressions can be enhanced?*

Secondly and very related to programming languages, if we take the Popularity of Programming Language (PYPL) Index<sup>2</sup> from June 2020 we can see that more than half of the share is occupied by languages that support the object oriented paradigm. And therefore this paradigm becomes the most used one. The aim of this paradigm is to model real world domains, according to [26]. That, in fact, is the same goal that ShEx has, it allows to model real world domains with schemas, and validate existing data with them. Therefore our second research relies on this and tries to automatically transform shape expressions into object domain models coded in any language that supports the object oriented paradigm:

**Research Question 2.** *Till which point can we automatically translate existing shape expressions in to object domain models?*

In order to answer the questions mentioned in this section, we will limit our working syntax to the reduced version of ShEx, defined in <sup>3</sup>. This reduced version is a strict subset of the full version and therefore, any conclusion drawn from it is extrapolated to ShEx.

## 1.2 Contributions

These are the major contributions of this dissertation:

1. A parser for the ShEx micro Compact Syntax. There are already existing parsers for ShEx and they work for ShEx micro Compact Syntax as it is a subset of ShEx, but they accept more structures than the ones defined by ShEx micro Compact Syntax. We propose a parser that is only focused on ShEx micro Compact Syntax and therefore error and warning messages can be enhanced.
2. Error and warning analyzer for schemas. Existing approaches do not semantically validate the schemas, they only perform error detection by means of complex grammars and parsers. Our proposed system does semantically validate the schemas by means of a custom analyzer that performs both syntactic and semantic analysis so it produces

---

<sup>2</sup><http://pypl.github.io/PYPL.html>

<sup>3</sup>[https://dcmi.github.io/dcap/shex\\_lite/micro-spec.html](https://dcmi.github.io/dcap/shex_lite/micro-spec.html)

human-friendly errors and warnings that users can use to fix their schemas.

3. Automatic translation of schemas in to object domain models in **Java** and **Python**. The proposed system integrates an open back-end with build-in code translation from the validated schemas to domain models in object oriented programming languages (*OOPL*) [5].
4. Evaluation of errors and warning generated of our proposed solution against existing tools. This comparison empirically shows the benefits and drawbacks of our proposed system.

### 1.3 Structure of the Document

The dissertation layout is as follows:

**Chapter 2** Indicates the state of the art of the existing RDF validation technologies, tools for processing Shape Expressions and other related projects.

**Chapter 3** Gives a basic theoretical background that it is needed to fully understand the concepts explained in the following chapters.

**Chapter 7** Contains a detailed initial planning and budget for the project, this is the designed planning followed during the execution of the project and the initial estimated budget.

**Chapter 8** Gives a basic theoretical background that it is needed to fully understand the concepts explained in the following chapters.

**Chapter 9** Provides a technical description of the design and implementation of the compiler itself. This includes, analysis, design, the technological stack choices, diagrams, implementation decisions and tests.

**Chapter 10** Compares the initial planning developed in chapter 4 with the final one. This includes the genuine execution planning of the project and the reasons and events that modified the one from chapter 4.

# CHAPTER 2

## Teoretical Background

For a proper understanding of this documentation and the ideas explained on it it is needed to know some theoretical concepts that are the fundaments of Linked Data, RDF, RDF Validation, programing languages and compilers. This sections is devoted to carefully explain those concepts to the needed depth to fully understand this dissertation, but for those readers that want a deeper explanation a more detailed view of the concepts presented here is offered in [19, 23, 7].

### 2.1 RDF

Resource Description Framework (RDF) is a standard model for data interchange on the web, started in 1998 and the first version of the specification was published in 2004 by the W3C according to [21]. RDF has features that facilitate data merging even if the underlying schemas differ, and it specifically supports the evolution of schemas over time without requiring all the data consumers to be changed. Another important feature is that RDF supports XML, N-Triples and Turtle syntax, the [Figure 2.1](#) shows an example of how a triplet can be written in RDF N-Triples Syntax.

RDF extends the linking structure of the Web to use URIs to name the relationship between things as well as the two ends of the link (this is usually referred to as a “triple” or "triplet"). Using this simple model, it allows structured and semi-structured data to be mixed, exposed, and shared across different applications. [2.3](#) shows an example of how different triples can be use to compose a graph, this graph represents the same as the [Figure 2.2](#)

This linking structure forms a directed, labeled graph, where the edges represent the named link between two resources, represented by the graph nodes. This graph view is the easiest

```
1 <http://example/subject1> <http://example/predicate1> <http://example/object1> .
```

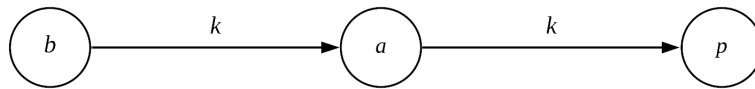
**Figure 2.1:** RDF N-Triples Example. From this example we can see that each triplet is composed of three elements, the subject the predicate and the object.

```

1 <http://example/bob> <http://example/knows> <http://example/alice> .
2 <http://example/alice> <http://example/knows> <http://example/peter> .

```

**Figure 2.2:** RDF N-Triples Graph Example. This exmaple shows the n-triples that generate the graph from [Figure 2.3](#).



**Figure 2.3:** RDF graph formed by triplets from [Figure 2.2](#), where *b* corresponds to `<http://example/bob>`, *a* corresponds to `<http://example/alice>`, *p* corresponds to `<http://example/peter>` and *k* corresponds to `<http://example/knows>`.

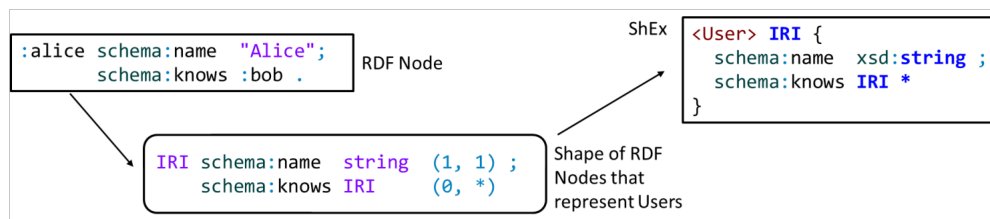
possible mental model for RDF and is often used in easy-to-understand visual explanations.

Also, related to this we strongly recommend the Tim Berners-Lee’s writings on Web Design Issues [9] where he explain the issues of the liked data and why is RDF so important.

## 2.2 Validating RDF

RDF therefore allows to represent and store data, and with this ability emerges the need to validate that the schema of the graph is correct. In order to perform the validation of RDF data there have been previous attempts, described in ??, this dissertation will focus on Shape Expressions. But in order to validate RDF data every technology will need to face the following RDF concepts:

- the form of a node (the mechanisms for doing this will be called “node constraints”);
- the number of possible arcs incoming/outgoing from a node; and
- the possible values associated with those arcs.



**Figure 2.4:** RDF node and its shape.

```

1 PREFIX :      <http://example.org/>
2 PREFIX schema: <http://schema.org/>
3 PREFIX xsd:    <http://www.w3.org/2001/XMLSchema#>
4
5 :User {
6   schema:name          xsd:string   ;
7   schema:birthDate     xsd:date?    ;
8   schema:gender        [ schema:Male schema:Female ] OR xsd:string ;
9   schema:knows          IRI @:User*
10 }

```

**Figure 2.5:** Shape Expression Example. This example describes a shape expression that describes a user as a node that has one name of type string, an optional birth date of type date, one gender of type Male, Female or free string and a set between 0 and infinite of other users represented by the knows property.

Figure 2.3 illustrates those RDF concepts by means of the Shape Expression that validates users. There we can see that the shape of the RDF node that represents Users represents the form of a node, the number of possible arcs and the possible value associated with those arcs.

### 2.2.1 Shape Expressions

As defined in [19] Shape Expressions (ShEx) is a schema language for describing RDF graphs structures. ShEx was originally developed in late 2013 to provide a human-readable syntax for OSLC Resource Shapes. It added disjunctions, so it was more expressive than Resource Shapes. Tokens in the language were adopted from Turtle and SPARQL with tokens for grouping, repetition and wildcards from regular expression and RelaxNG Compact Syntax [25]. The language was described in a paper [23] and codified in a June 2014 W3C member submission which included a primer and a semantics specification. This was later deemed “ShEx 1.0”.

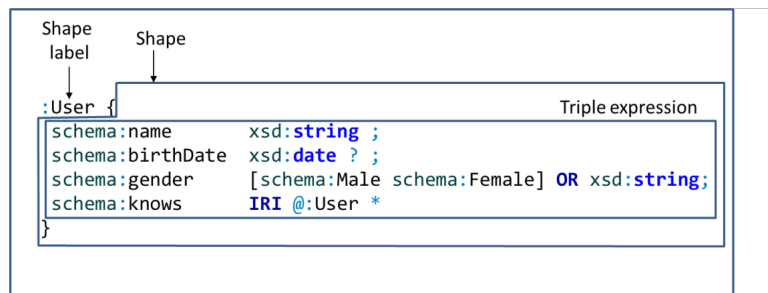
As of publication, the ShEx Community Group was starting work on ShEx 2.1 to add features like value comparison and unique keys. See the ShEx Homepage <http://shex.io/> for the state of the art in ShEx. A collection of ShEx schemas has also been started at <https://github.com/shexSpec/schemas>.

#### ShEx Compact Syntax: ShExC

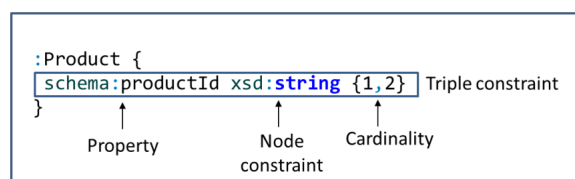
The ShEx compact syntax (ShExC) was designed to be read and edited by humans. It follows some conventions which are similar to Turtle or SPARQL.

- PREFIX and BASE declarations follow the same convention as in Turtle. In the rest of this chapter we will omit prefix declarations for brevity.
- Comments start with a # and continue until the end of line.





**Figure 2.6:** Shapes, shape expression labels and triple expressions.



**Figure 2.7:** Parts of a triple expression.

- The keyword `a` identifies the `rdf:type` property.
- Relative and absolute IRIs are enclosed by `< >` and prefixed names (a shorter way to write out IRIs) are written with prefix followed by a colon.
- Blank nodes are identified using `_:label` notation.
- Literals can be enclosed by the same quotation conventions ( `'`, `"`, `'''`, `"""` ) as in Turtle.
- Keywords (apart from `a`) are not case sensitive. Which means that `MinInclusive` is the same as `MININCLUSIVE`.

A ShExC document declares a ShEx schema. A ShEx schema is a set of labeled shape expressions which are composed of node constraints and shapes. These constrain the permissible values or graph structure around a node in an RDF graph. When we are considering a specific node, we call that node the focus node.

Figure 2.6 shows the first level of a shape expression, we have a label and the shape itself that is what we asing to the `:User` label. Then, the shape is composed by triple expressions. The triple expression structure is explained in ??, and as its name indicates it is composed of three elements, the property, the node constraint and the cardinality.

Shape Expressions Compact Syntax is much bigger and contains other multiple features that give ShEx its power, and all of them can be explored in [19] but they are not needed to understand this dissertation.

## Use of ShEx

Strictly speaking, a ShEx schema defines a set of graphs. This can be used for many purposes, including communicating data structures associated with some process or interface, generating or validating data, or driving user interface generation and navigation. At the core of all of these use cases is the notion of conformance with schema. Even one is using ShEx to create forms, the goal is to accept and present data which is valid with respect to a schema. ShEx has several serialization formats:

- a concise, human-readable compact syntax (ShExC);
- a JSON-LD syntax (ShExJ) which serves as an abstract syntax; and
- an RDF representation (ShExR) derived from the JSON-LD syntax.

These are all isomorphic and most implementations can map from one to another. Tools that derive schemas by inspection or translate them from other schema languages typically generate ShExJ. Interactions with users, e.g., in specifications are almost always in the compact syntax ShExC. As a practical example, in HL7 FHIR, ShExJ schemas are automatically generated from other formats, and presented to the end user using compact syntax.

ShExR allows to use RDF tools to manage schemas, e.g., doing a SPARQL query to find out whether an organization is using `dc:creator` with a string, a `foaf:Person`, or even whether an organization is consistent about it.

## ShEx Implementations

[Check links.](#)

At the time of this writing, we are aware of the following implementations of ShEx.

- `shex.js` for Javascript/N3.js (Eric Prud'hommeaux) <https://github.com/shexSpec/shex.js/>;
- `Shaclex` for Scala/Jena (Jose Emilio Labra Gayo) <https://github.com/labra/shaclex/>;
- `shex.rb` for Ruby/RDF.rb (Gregg Kellogg) <https://github.com/ruby-rdf/shex>;
- Java ShEx for Java/Jena (Iovka Boneva/University of Lille) <https://gforge.inria.fr/projects/shex-impl/>; and
- `ShExkell` for Haskell (Sergio Iván Franco and Weso Research Group) <https://github.com/weso/shexkell>.

There are also several online demos and tools that can be used to experiment with ShEx.

- `shex.js` (<http://rawgit.com/shexSpec/shex.js/master/doc/shex-simple.html>);
- `Shaclex` (<http://shaclex.herokuapp.com>); and
- `ShExValidata` (for ShEx 1.0) (<https://www.w3.org/2015/03/ShExValidata/>).

### 2.2.2 Other Technologies

As other validation technologies we will just explore the existence of them as it is very interesting to know how other tools approach the same issue.

#### SHACL

Also in [19], Chapter 5, it is fully explained that Shapes Constraint Language (SHACL) has been developed by the W3C RDF Data Shapes Working Group, which was chartered in 2014 with the goal to “produce a language for defining structural constraints on RDF graphs [24].”

The main difference that made us choose ShEx over SHACL are that ShEx emphasized human readability, with a compact grammar that follows traditional language design principles and a compact syntax evolved from Turtle.

#### JSON Schema

JSON Schema born as a way to validate JSON-LD, and as turtle and RDF can be serialized as JSON-LD it is usual to think that JSON Schema can validate RDF data, but this is not fully correct. And the reason is that the serialization of RDF data in to JSON-LD is not deterministic, that means that a single schema might have multiple serializations, which interferes with the validation as you cannot define a relative schema.

## 2.3 Programming Languages

According to [7] “a programming language is a formal language comprising a set of instructions that produce various kinds of output.” When we talk about programming languages we need to know that they are split into two, General Purpose Languages (GPL) and Domain Specific Languages (DSL). The main difference overtime is that, as said in [2], a domain-specific language (DSL) is a computer language specialized to a particular application domain in contrast to a general-purpose language (GPL), which is broadly applicable across domains.

In the specific case of ShEx-Lite we will be talking about a Domain Specific Language, and more deep we would classified it as a Declarative one, that means that it is not Touring Complete [8].

## 2.4 Compilers

A compiler is a computer program that translates computer code written in one programming language (the source language) into another language (the target language). Is during this translation process where the compiler validates the syntax and the semantics of the program, if any error is detected in the process the compiler raises an exception (understand as a compiler event that avoids the compiler to continue its execution).

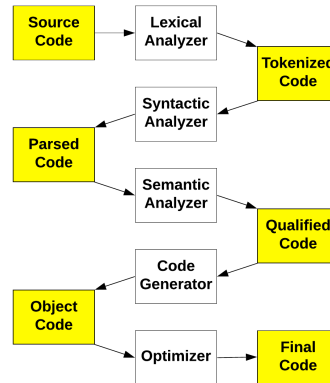


Figure 2.8: Compiler stages.

### 2.4.1 Internal Structure

In order to decompose the internal structure of a compiler they have been split in to the most common task they do [Figure 2.8](#), of course this doesn't mean that there are compilers with more or less stages, but at the end everything can be group into any of the groups that we will explain:

#### Lexical Analyzer

The lexical analyzer task is to get the input and split it in to tokens [4], which are build from lexemes. If the compiler cannot find a valid token for some lexemes in the source code will generate an error, as the input cannot be recognized.

#### Syntactic Analyzer

The syntactic analyzer takes the tokens generated during the lexical analysis and parses them in such a way that try's to group tokens so the conform to the language grammar rules. During this stage if there is any error while trying to group the tokens then the compiler will rise an error as the input cannot be parsed.

#### Semantic Analyzer

The semantic analyzer has two main tasks, usually. First it validates that the source code semantics are correct, for example `4 + "aaa"` would not make sense. And the second task is to transform the Abstract Syntax Tree in to a type-checked and annotated AST. Usually that means relate the invocations and variables to its definition, very useful for type-checking.

### Code Generator

The task of the code generator as its name indicates is to generate the target code, it can be byte code, machine code or even another high-language code.

### Code Optimizer

The code optimizer is the last step before the final target code is generated, it rewrites the code that the code generator produced without changing the semantics of the program, its aim is just to make code faster. At [6] you can see an example of some optimizations that can be done at compile time to make your code faster.

#### 2.4.2 Conventional Compilers

Conventional compiler are a big monolith where each stage 2.8 is executed automatically after the previous stage, if the compiler has eight steps you need to execute them all at once. This approach have been the “old-fashion” but it presents some drawbacks:

- A poor IDE [3] integration. IDE’s need to perform incremental compilations in matter of nanoseconds so the user doesn’t feel lag when typing the program. With conventional compilers as you need to go through all the compilation process at once they where very slow and companies like Microsoft need to develop different compilers, one for the IDE and another for the final compilation of the program itself. This lead to several problems like that if a feature gets implemented in the final compilation compiler but not in the IDE one the IDE would not support the feature meanwhile the language would.
- Difficult to debug. As the conventional compilers where a blackbox the only way to test intermediate stages was by throwing an input and waiting the the feature you wanted to test was thrown for that input.

#### 2.4.3 Modern Compilers

After the problems Microsoft had with the C# compiler they decide to rewrite the whole compiler and introduce a concept called “compiler as an API” with Roslyn [1]. This concept has been perfectly accepted and solved many problems. In this concept each stage has an input and an output that can be accessed from outside the compiler and stages can be executed independently on demand. This means that for example if an IDE just want to execute the Lexer the Parser and the Semantic analysis it can. That translates in to speed for the user.

Also the second problem is solved as testing individual parts of the compiler is much more easy than the hole compiler at once.

# CHAPTER 3

## Related Work

Some work has already been done in the field of Shape Expressions and RDF validation technologies. In this chapter we will go over the main studies related to our project, exploring what they have achieved and some of their limitations.

### 3.1 Simplifications of ShEx

#### 3.1.1 The S language

In 2019 at [18] was defined a language called **S** as a simple abstract language that captures both the essence of ShEx and SHACL. This is very relevant as this language is intended to be the input of a theoretical abstract machine that will be used for graph validation for both ShEx and SHACL. Also in the same paper the authors carefully describe the algorithm for the translation from ShEx to S and from SHACL to S.

Although the theoretical abstract machine has not been implemented yet the intention of the WESO Research Group, where this S language was defined, is to devote more efforts in to this project during the 2021.

Other definition of an abstract language based on uniform schemas can be found at [11]. This language is focused on schemas inference rather on validation, but needs to be taken into account as they also perform an abstraction of both ShEx and SHACL.

#### 3.1.2 ShExJ Micro Spec

Recently the Dublin Core Team<sup>1</sup> is working into an specification that allows to define Shape Expressions in tabular formats. For this specification they propose a simplification of the Shape Expressions JSON syntax that allows to define an schema as a set of simple triple constraints. This specification is not official and has not been validated yet but it is very important for our work as we will also work in a simplification of a syntax of ShEx.

And to the best of our knowledge and after the research process carried out for this project no

---

<sup>1</sup><https://dublincore.org/>

other language based on a subset of Shape Expressions has been designed nor implemented yet.

## 3.2 ShEx Ecosystem Tools

We already know that ShEx and SHACL have been the two main technologies for RDF validation and some tools emerged around them, we think that some of them might benefit from ShEx-Lite. Here we introduce briefly those that had the biggest impact in the community.

### 3.2.1 Validators

Since the beginning of ShEx and SHACL as languages the RDF community started to build tools that take as input the schemas defined and validate graphs.

This kind of tools can benefit from ShEx-Lite from the point of view that new functionalities can be easily implemented and tested in the lite version of the language before even touching the stable releases of both tools. In the case of ShEx this is more obvious as ShEx-Lite and ShEx are both implemented in Scala and if good design principles are used functionalities can be just migrated and expanded for the rest of the language.

The most important validators are:

#### Shaclex

According to the Shaclex<sup>2</sup> official website it is an Open Source Scala pure functional implementation of an RDF Validator that supports both Shape Expressions and SHACL. It was initially developed by Dr. Jose Emilio Labra Gayo and is being maintained by an active community on GitHub. It is used by different projects around the globe and its goal is to validate RDF graphs against schemas defined in Shape Expression or in SHACL.

This implementation of a ShEx validator is very important for us as ShEx-Lite is completely inspired by it and aims to transfer the syntactic and semantic validation enhancements to it.

#### ShEx.js

Another example of as a ShEx validator implementation is `ShEx.js` which is JavaScript based and also open source on GitHub. This implementation is very important for the ShEx community as they defined the serialization of the AST in this implementation as the abstract syntax of ShEx.

---

<sup>2</sup><https://github.com/weso/shaclex>

### 3.2.2 IDEs

In order to facilitate the task of writing schemas some engineers decide to implement specific IDEs for the Shape Expressions Language.

This tools will completely benefit from ShEx-Lite and there are currently collaborations in process. At the time they work with Shaclex, which is structured as a conventional compiler, but with the API architecture of ShEx-Lite IDEs can access directly to the syntactic and semantic modules so features like advances coloring syntax or incremental compilation are available.

#### YASHE

YASHE<sup>3</sup> (Yet Another ShEx Editor), is a Shape Expressions IDE which started as a fork of YASQE(which is based on SPARQL). This tool performs lexical and syntactic analysis of the content of the editor, thus offering the user a realtime syntactic error detector. It has features like: syntax highlighting, visual aid elements (tooltips) and autocomplete mechanisms. In addition, it offers a simple way of integrating into other projects.

#### Protégé

Protégé is a piece of software developed by the University of Stanford focused on ontology edition. During the last year they added support for Shape Expressions dition on their own software so they became another ShEx IDE.

#### VSCode

VSCode is a source code light-weight editor developed by Micorsoft and supported by Linux, macOS and Windows. By default this editor does not support any programming language, the way it works is with packages that the community develops and extends the functionality. One of those packages adds support for Shape Expressions Compact syntax and transforms VSCode into a ShEx IDE.

This plugin does not add semantic validation and it is a clear target to benefit from ShEx-Lite features.

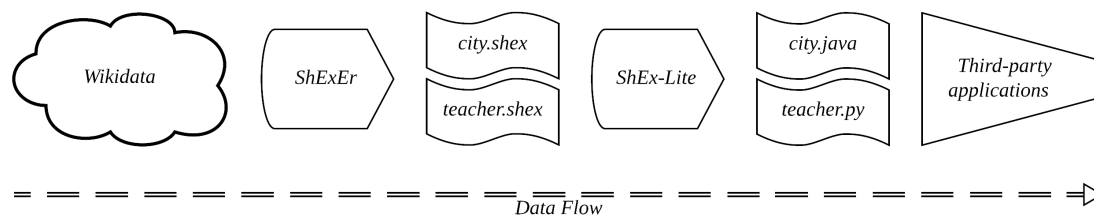
### 3.2.3 Others

Other researches focused their efforts in to inferring schemas to existing data sets and creating tools to that evolved from ShEx in order to transform existing data.

---

<sup>3</sup><https://github.com/weso/YASHE>





**Figure 3.1:** ShEx-Lite integration with Shexer for automatically generating java domain object models for the Wikidata schemaless existing data. This shows the schemaless data from wikidata from which shape expressions are inferred by shexer and later transformed to java plain objects by means of ShEx-Lite so third party applications can implement the domain model.

### Shexer

Shexer<sup>4</sup> is a python library aimed to perform automatic extraction of schemas in both ShEx and SHACL from an RDF input graph. That is if all the other tools take the schemas as the input and validate a graph with it, this tool takes a graph and from it it infers the schemas that it might contain. Its work is fully described in [11, 13].

### ShExML

ShExML<sup>5</sup> is a language based on ShEx (not a simplification nor an abstraction of ShEx) that can map and merge heterogeneous data formats into a single RDF representation. The main idea behind this tool is written at [16].

An example of how this different tools can work together thanks to ShEx-Lite would be the following, illustrated at Figure 3.1. Wikidata currently holds millions of registers that do not have any schema that validates them. And they need to make consumer that represents the data in to an object domain model. Without any tool this is just almost impossible, but this shexer you can infer the schemas to ShEx-Lite syntax and with the ShEx-Lite compiler you can automatically create the object domain model in your favorite OOL.

<sup>4</sup><https://github.com/DaniFdezAlvarez/shexer>

<sup>5</sup><https://github.com/herminiogg/ShExML>

## Part I

# Enhancing Error and Warning Detection and Emission on ShEx

# CHAPTER 4

## Analysis of Existing Sintactic and Semantic Analizers

In the Related Work ([Chapter 3](#)) some ShEx tools were explained. This section will detail more those tools that provide any kind error and warning detection and emition. After, we will detail the points that we think can be enhanced.

Before start the analysis we must define a methodology in order to be able to make an even analysis for all existing tools.

### 4.1 Methodology

To evaluate existing systems from a neutral point of view we will use the ShEx specification as the basis. However, this specification does not cover all possible cases, in particular it leaves most semantic restrictions to the choice of the specific implementation.

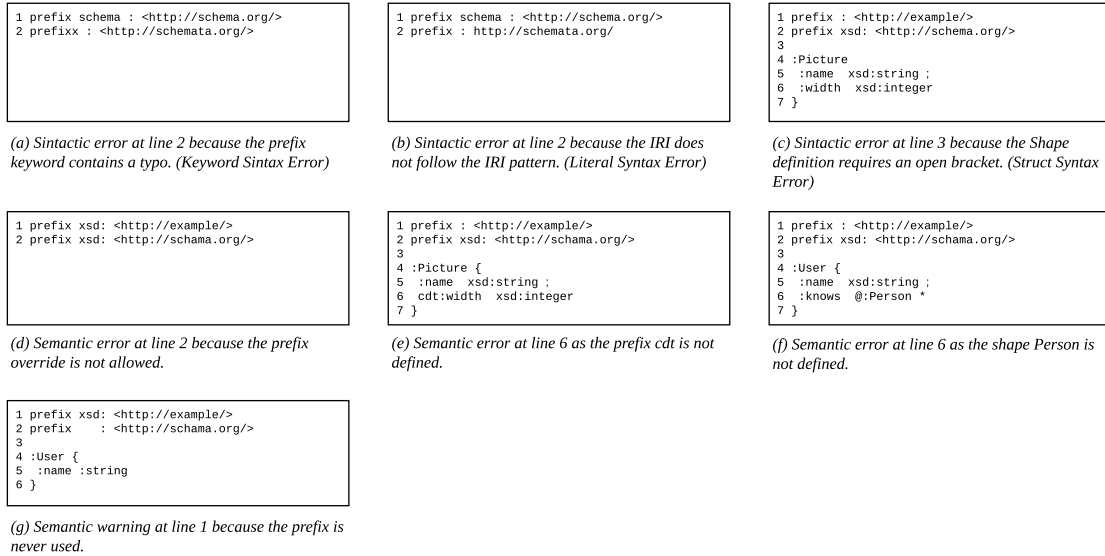
Therefore, as regards this evaluation, when a semantic option not contemplated by the specification is proposed, the option that favors the security of the language will be chosen. For example. If the specification did not say anything about whether a variable can be redefined and we had to take an option, we will always choose not, so that the language is as safe as possible and does not lead to errors.

The unique sintactic restrictions applied is:

- In the last triple constraint of a set expression the trailing semicolon it is optional but recommended.

The semantic restrictions that have been applied are listed below.

- Overwriting of prefixes is not allowed.
- Overwriting of the base is not allowed.
- Overwriting of the start shape is not allowed.
- Overwriting of shapes is not allowed.



**Figure 4.1:** Examples of ShEx micro Compact Syntax code containing errors.

- All references must exist within the scope of the schema.

In addition, in this evaluation we will use different test cases for each system, specifically the test cases correspond to each element of the ShEx micro Compact grammar. Remember that the elements that this grammar has are: *definition of prefixes*, *definition of the base*, *definition of the start shape* and *definition of shapes*. To others within the previous elements you will also find references to prefixes, the base and other shapes. Therefore we will test all these elements in their syntactic and semantic aspects. Figure 4.1 shows some examples of this errors.

## 4.2 Sintactic Analyzers

According to [14] we consider a Sintactic Analyzer a piece of software capable of parse, generate a parse tree and detect and emmit sintactic warnings and errors.

Therefore in this category we would include **Shaclex**, **ShEx.js**, **YASHE** and **VS Code Plugin**. Table 4.1 shows a comparison between the analyzed tools.

Some comments to be made about the results obtained are that although we get an error for syntactic errors, the quality of the error is more or less always the same. For example for the fragment `prefixx xsd: <http://example/>` where we introduced an error at the keywork `prefix` by adding an extra `x` the error obtained is: `This line is invalid. Expected: PNAME_NS.`

To our point of view this error message nor is not correct because it does not provide the

**Table 4.1:** Detection of the different syntactic errors by the current existing ShEx tools that syntactically analyze the shape expressions.

Syntactic Errors								
Analyzers	Prefix Definition	Base Definition	Start Shape	Shape Definition	Prefix Reference	Base Reference	Shape Reference	Recommends Semicolon Last Triple Constraint
Shaclex	Yes	Yes	Yes	Yes	Not completely	Yes	Yes	No
ShEx.js	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
YASHE	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No
VS Code Plugin	Yes	Yes	Yes	Yes	Yes	Yes	Yes	No

**Table 4.2:** Detection of the different semantic errors by the current existing ShEx tools that semantically analyze the shape expressions.

Semantic Errors								
Analyzers	Prefix Override	Base Override	Start Shape Override	Shape Override	Non Existing	Non Existing	Non Existing	
					Prefix Reference	Base Reference	Shape Reference	
Shaclex	No	No	No	No	Yes	-	No	
ShEx.js	No	No	No	No	Yes	-	No	
YASHE	No	No	No	No	Yes	-	Yes <sup>1</sup>	

user enough information to fix the schema.

Then also it is important to remark that during this analysis we encounter other syntactic problems that were not detected by tools like Shaclex, an example is that properties like `schema:rdf@:name` (which is not a valid IRI) are accepted without errors.

### 4.3 Semantic Analyzers

As Semantic Analyzers we will only consider those tools that validate the semantics of the language, in this section we include the validation of references like prefixes and shapes. The tools that claim to support this validations are **Shaclex**, **ShEx.js**, and **YASHE**. Table 4.2 shows a comparison between the analyzed tools.

From the obtained results we have to point that most of the tools opted for an open policy when talking about language semantics. From our point of view this has its advantages and its drawbacks. But this only affects to the override policy. All of the tools should implement the non existing references validation and most of them only focus on prefixes definition with the exception of YASHE which does the checking of the shape reference but the error message sometimes is not completely accurate.

It is also remarkable that none of the tools performs a deeper analysis so there is no detection of unused resources, therefore no warnings are generated by none of the existing tools.

## 4.4 Possible Enhancements

Previous sections show the current state of the existing tools, their capabilities and their lacks. With all that information we propose a list of enhancements that can be done to improve the error and warning detection. As seen in previous sections there's work that can be done to improve the existing ecosystem of tools. We have identified the following aspects that will benefit end users:

- **Enhancement of error messages.** Existing error messages, originated both by syntactic or semantic errors do not offer information about the exact place that originates the error nor a processed description nor possible solutions.
- **Creation of a new type of error messages with lower importance called warnings.** Currently systems do not analyze if declared resources are used and therefore there is no need to generate warnings. We propose to not only fully analyze the resources to detect non-used ones but also the creation of error messages with lower importance like warnings that can be used to offer more information to the end user.
- **Detection of override definitions.** Most of the existing tools prefer not to detect when a definition is being overridden, we propose to detect those situations and treat definitions as fixed values.
- **Detection of undefined references.** Some tools detect some broken references, we propose to enhance this situation and take that behaviour to other elements like shape references.
- **Detection of unused resources.** Related to the second point sometimes new users copy and paste old code which ends with lots of unused code, we propose a system that detects those situations and suggest to remove that unused code.

# CHAPTER 5

## Proposed Sintactic and Semantic Analyzer

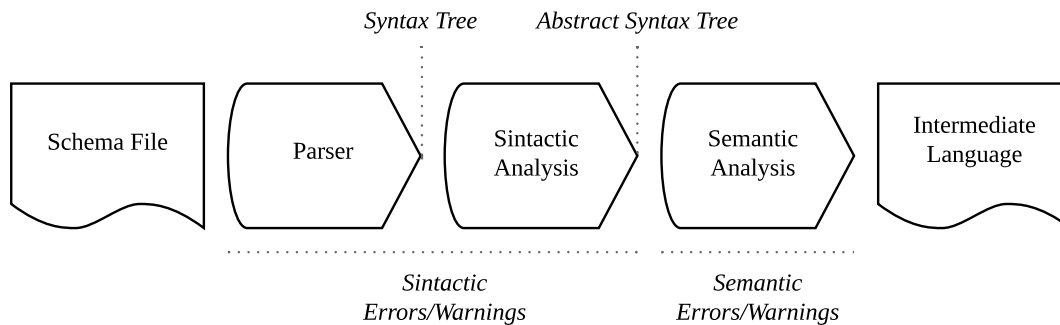
Once all the objectives and requirements to be achieved have been described, the different systems and techniques existing to achieve them have been studied, and their contributions and shortcomings have been evaluated, we will describe the proposed solution both in terms of design and possible implementation

### 5.1 Structure

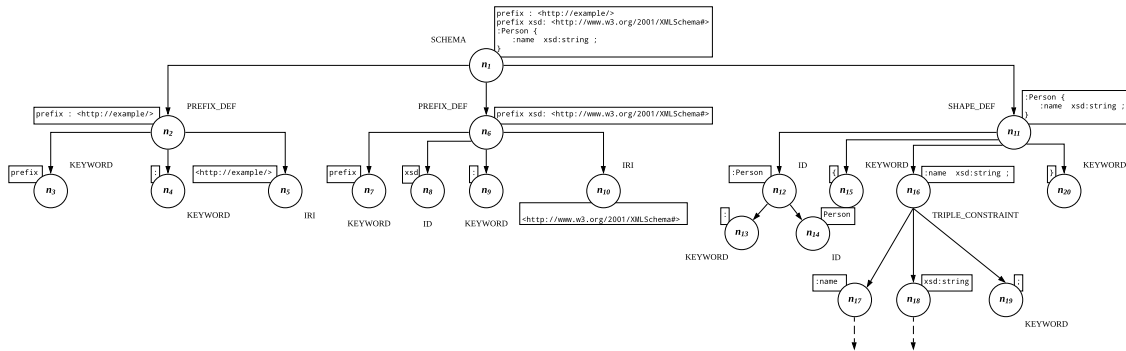
The system is divided into components so that each component works on its input and produces its output. In this way, a parser is achieved that behaves like an API where each element can be called individually. [Figure 5.1](#) shows the different components of this analyzer.

#### 5.1.1 Parser

We define the parsing stage as the process that begins when we receive the source code that makes up the schema until the moment we produce a syntax tree. Therefore it includes the conversion to tokens by the lexer, the grouping of tokens in rules and later in a syntax tree



**Figure 5.1:** Sintactic and Semantic Analyzer structure.



**Figure 5.2:** Syntax Tree twenty first nodes produced by the parser.

by the parser.

The general idea of this stage is that you take the source code as input and build a syntactic tree with all the possible information from the source code. This implies that the syntactic tree is not only made up of abstract grammar, but also of separators, braces and keywords. [Figure 5.2](#) shows an example of the first 20 nodes generated by the parser. There we can see this composition of separators, keywords, braces and content.

Once we have the complete syntactic tree generated, we can go through it to carry out syntactic analysis on the different elements. For example, in the tree in [Figure 5.2](#) we could implement a validator that in the event that the last triple constraint of a shape definition (*node 16*) did not have the semicolon termination keyword (*node 19*), it would generate a warning message to the user.

### 5.1.2 Sintactic Analyzer

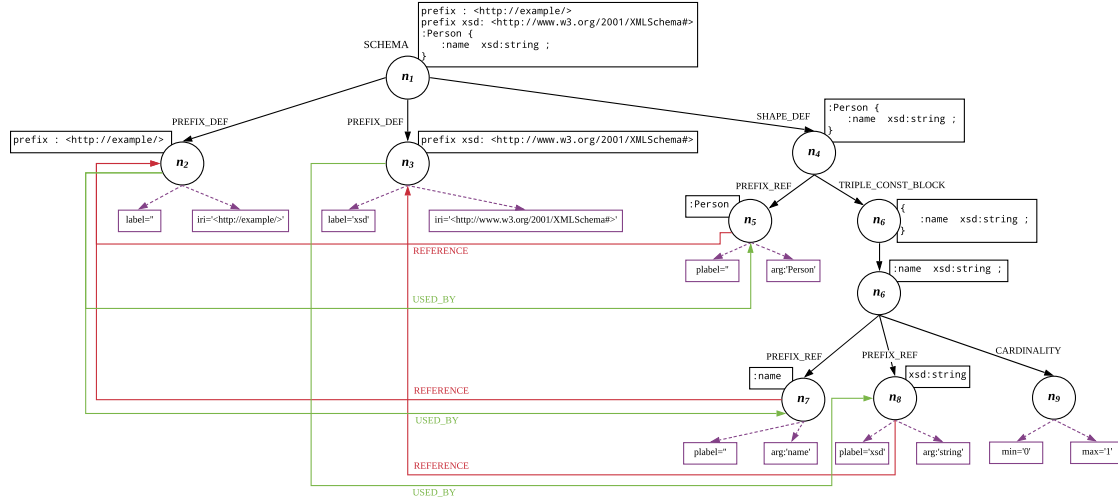
The sintactic analyzer is in charge of traversing the syntactic tree in order to search for possible patterns that the user has to be informed about. If none were found it would be understood that the syntactic tree is well formed and it will tranform the Syntax Tree [Figure 5.2](#) into an Abstract Syntax Tree [Figure 5.3](#) (*without the green and red relations*).

For this, each node within our syntactic tree is aware of the context in which it is. Therefore we can ask questions to the nodes, such as to a prefix definition ([Figure 5.2](#) *n<sub>2</sub>*), do you have a label? (*No*) or who is the node that defines your iri? (*node n<sub>5</sub>*). With questions like these, the syntactic tree can be analyzed for patterns that represent warnings or errors.

### 5.1.3 Semantic Analyzer

The semantic analyzer is responsible of building all the possible relations between the AST nodes, analyze and check that all those relations that must exist indeed exist. For this porpouse as just seen we reduce our Syntax Tree to an Abstract Syntax Tree. [Figure 5.3](#)





**Figure 5.3:** Abstract Syntax Tree produced after validation and transformations.

Shows a the resulting AST after the corresponding analysis and transformations, we call this graph the *Intermediate Language*.

Once we have the representation modeled and this representation is capable of expressing all the assumptions of our language, we can begin to apply validators on our structure. For example if we wanted to find broken references we could go to the nodes that are a reference to definitions like nodes  $n_5, 7$  and  $8$  and check that there is indeed a valid reference for each of them.

Furthermore, we can even analyze how many times a definition is used by a reference so that we can launch messages warning the end user in some cases, such as when a prefix is not used.

## 5.2 Implementation

As proof of concept of the previous proposal we offer an implementation of the three components, the parser, the syntactic analyzer and the semantic parser. The implementation is defined in the same way as the structure, in three parts. We will now explore each of those parts and their responsibilities separately.

### 5.2.1 Parser

As previously discussed, the function of the parser is to extract a syntactic tree from the diagrams that we can analyze. For this purpose we decided to use the Antlr tool [22]. This tool is capable of generating syntactic analyzers from grammars defined in its own syntax. However, this tool is focused on completely processing the syntax tree and producing only

```
1  override def visitConstraint_triple_expr(...) {  
2      if(/*No trailing semicolon*/)  
3          //Warn user about this bad practice  
4  }
```

**Figure 5.4:** Checker implementation for missing semicolons warning generation.

the abstract syntax tree. Therefore we had to use a modification of the original ShEx micro Compact Syntax syntax so that Antlr would produce a tree with all the syntactic content. This also does offer the flexibility that in the future if we want to implement any additional syntactic validation we simply have to do it on the tree that the parser generates for us and not on the Antlr code.

### 5.2.2 Sintactic Analyzer

The sintactic analyzer has the responsibility to validate that the parser produced syntax tree is correct and to build the abstract syntax tree as well. To do this, it uses the same mechanism. Through the visitor pattern we go through our syntax tree. Each implementation of this visitor has a purpose, for example an implementation can go through a few specific nodes to validate them syntactically while another can go through them in order to build the AST. [Figure 5.4](#) shows an example of how a sintactic check is implemented.

The AST construction stage is very delicate since for each generated node we have to include as much context information as possible so that when an error is detected in the tree we can identify not only the cause but also the position, the origin, the rest of the affected nodes and therefore offer a content-rich error message. Regarding our implementation, for each node we save the following context information:

- **Source file path.** Represents the path to the source file where the node was generated.
- **Line.** The line in the source file where the node was generated.
- **Column.** The column in the source file where the node was generated.
- **Token interval.** The interval (*start, end position*) of tokens from the source file that generated the node.
- **Content.** The content of the node as plain text. [Figure 5.3](#) is very representative of this.
- **Parent node.** A pointer to the parent node.
- **Children nodes.** A list of pointers to all the children nodes.

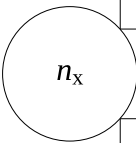
[Figure 5.6](#) represents this information inside each node. Our default implementation only looks for the following extra syntactic pattern to the other implementations seen in [Chapter 4](#):

```

1 warning[W005]: missing semicolon
2 --> ../../test/assets/input_correct_java_schema_user_car_1.shexl:8:23
3 |
4 8 | :knows      @:User *
5 |           ~ semicolons are not compulsory in the last triple constraint,
6             but its usage its encouraged as otherwise your code wont be
7             following shape expressions specification.

```

**Figure 5.5:** Syntactic warning produced by the proposed syntactic analyzer.



<i>Property</i>	<i>Value Type</i>
<i>Source File</i>	<i>String</i>
<i>Line</i>	<i>Integer</i>
<i>Column</i>	<i>Integer</i>
<i>Token Interval</i>	<i>Interval&lt;Integer&gt;</i>
<i>Content</i>	<i>String</i>
<i>Father</i>	<i>Node</i>
<i>Children</i>	<i>List&lt;Node&gt;</i>

**Figure 5.6:** Common information stored at any AST node.

shape expressions whose last triple constraint does not contain the semicolon ending character. In case we find this pattern, we inform the event manager that a notice has been found that must be passed on to the user, [Figure 5.5](#).

### 5.2.3 Semantic Analyzer

Recall that the semantic analyzer takes the generated AST, runs it in search of errors and transforms it in such a way that it emits a graph that corresponds to the intermediate language. We can separate semantic analysis into two phases, a first one in which we transform our syntactic tree, adding possible relationships. And a second phase in which we analyze existing and created relationships.

#### Tree transformations

In the case of our syntax the semantic relations that we find is the linking of a reference to its definition and the opposite direction to indicate that a definition is being referenced by a node. The transformations are listed bellow:

- **Linking prefix definition with prefix references.** Prefix references occurs when a node describes itself as the composition of a prefix and an argument. The idea is

that the prefix substitute the IRI, but must be linked as any prefix reference needs to point to an existing definition.

- **Linking base definition with base references.** Some nodes are defined as relative IRIs to the base definition and therefore need to be linking to them in order to be able to get that base IRI.
- **Linking shape definition with a shape reference.** Shape definitions can be used at the `start` definition to point the default shape or as type constraints in the triple constraints. At any of those points shape references must exist within the scope of an schema.

### Tree relations analysis

For this purpose, the semantic analyzer defines the visitor pattern on the nodes of the abstract syntax tree so that each of the different analysis is done with a tree visiting implementation. Some of the

## 5.3 Sintactic and Semantic Error and Warnings Detected

### 5.3.1 Not trailing semicolon at last triple constraint

### 5.3.2 Prefix not defined

### 5.3.3 Shape not defined

### 5.3.4 Prefix overridden

### 5.3.5 Shape overridden

### 5.3.6 Unused prefix definition

### 5.3.7 Base set but not used

## Part II

# Translating ShEx Schemas to Object Domain Models

# CHAPTER 6

## Object Domain Model Translation Problem

The ODMTP (*Object Domain Model Translation Problem*), when talking about Shape Expressions, is the aim to transform existing schemas, that already represent domain models, in to object domain models. Or what it is the same, translate the ShEx schemas to objects coded in some Object Oriented Language. [Figure 6.1](#) represents this aim. The problem is to convert the *Source* in to the *Target*.

Person Schema (Source)	Person Java Object (Target)
1 # Prefixes...	1 // Imports...
2 :Person {	2 public class Person {
3     :name xsd:string ;	3     private String name;
4     :knows @:Person *	4     private List<Person> knows;
5 }	5     // Constructor...
	6     // Getters and Setters...
	7 }

**Figure 6.1:** Schema modeling a **Person** in ShExC syntax to the left. And the expected translated code in Java to the right.

This problem, with the example from [Figure 6.1](#) might seem simple to solve, but before proposing a solution we will need to explore if the all that ShEx can express can be translated to every object oriented programming language. For that propose we will generalize the generated domain models to what it is known as PO (*Plain Objects*) [15], that is, simple objects that do not rely on any framework. And then the main question to this problem changes to if we can translate any Shape Expression to a Plain Object.

In order to solve that question we will need to explore and compare the expressivity of both concetps.

### 6.1 Shape Expressions Expressivity

Best way to study Shape Expressions expressivity is by exploring its data model and everything that can be represented by them, [Figure 6.2](#) displays the ShEx recursive data model.

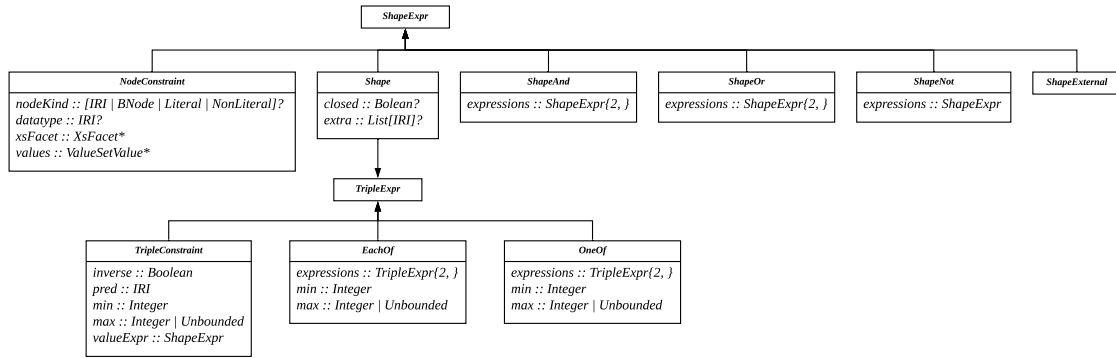
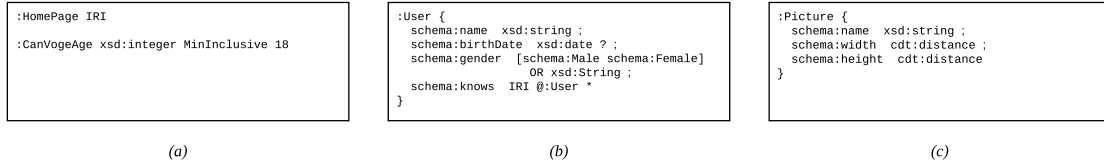


Figure 6.2: ShEx data model.



**Figure 6.3:** ShEx example schemas. (a) applies a node constraint to a root node. (b) Shows how an conditional constaint can be applied to a property node. And (c) shows how custom data types ctd can be applied to a property node.

As we can see from [Figure 6.2](#) the data model of ShEx is hardly complicated, it contains, conditional constraints, root node constraints, nested constraints and even logical constraints, [Figure 6.3](#) shows some examples of this.

When we talk about translating schemas into object domain models we are only interested in the subset of shape expressions that are formed by a set of properties, that means the shape expressions that define constraints for root nodes, like [Figure 6.3 \(a\)](#), won't be of our interest.

### 6.1.1 Properties Set Defined Shapes

A property set defined shape is a shape expression that it is defined as a set of properties, for example [Figure 6.3 \(b\)](#) and (c). This kind of shape expressions express that a node must contain the indicated properties, each one with its type.

Then, we can generalize that a shape expression is formed by the schema name and the list of constraint-annotated properties, [Figure 6.4](#) shows the different parts that compose a property set defined shape.

Furthermore, we can conclude:

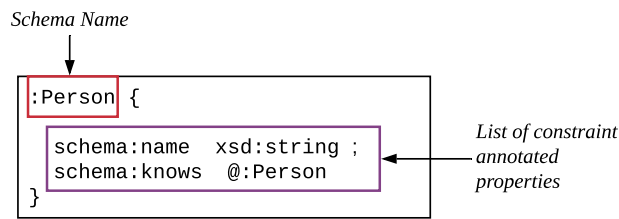


Figure 6.4: Property set defined shape model.

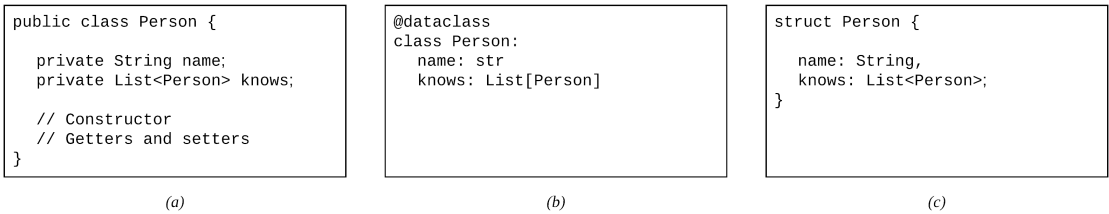


Figure 6.5: Java, Python and Rust codings of Person object. *a* corresponds to Java, *b* corresponds to Python and *c* corresponds to Rust.

**Partial Conclusion 6.1.** *The expressivity of a property set defined shape depends on the node constraints of its properties.*

6.1.2 Property Node Constraints

The different constraints that a node might have in Shape Expressions are shown in .

Insertar tab  
node constr  
del libro de

6.2 Plain Objects Expressivity

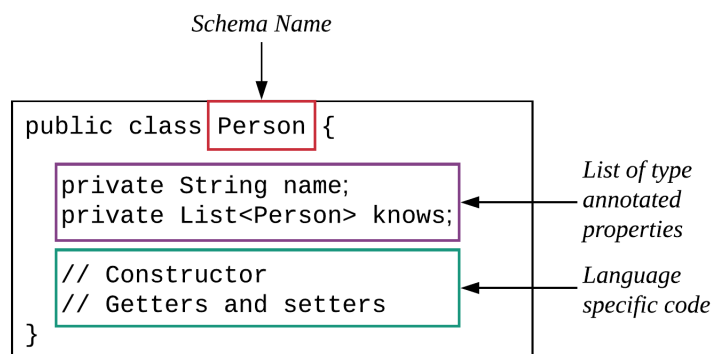
Plain objects can be coded in any object oriented programming language, or at least in any language that supports this paradigm. First we will explore how plain objects are generally coded, then how the language increases or decreases the expressivity and finally we will generalize the core concepts that can be expressed by any plain object codification.

6.2.1 Plain Objects Structure

From the existing programming languages we can infer the general structure of plain objects. For this porpuse we take the PYPL Index (*PopularitY of Programming Language*)<sup>1</sup> from June 2020 and take the 2 most used programming languages that support the object oriented

<sup>1</sup><http://pypl.github.io/PYPL.html>





**Figure 6.6:** Java plain object decomposition.

paradigm, those would be Java and Python. And then, just to enlarge the scope we will take Rust because it is a new programming language that includes lots of features.

Figure 6.5 shows three models that correspond to the codification of the Person schema from Figure 6.1. For example if we analyze the Java fragment, that seems to be the most complex one out of the three fragments we can see in Figure 6.6 that it is composed by the *Schema Name*, the *List of Type Annotated Properties* and some *Language Specific Code*. This correlates to the other two programming languages as they also contain this three elements. Therefore after this brief analysis we can conclude:

**Partial Conclusion 6.2.** *Any Plain Object, coded in a programming language that supports the object oriented paradigm, will be composed of those three elements: **Schema Name**, **Type Annotated Properties** and **Language Specific Code**.*

### 6.2.2 Plain Objects Language Expressivity Dependence

In Figure 6.5 we can see that all of three languages use similar types to represent the Person model. But with just one example we cannot generalize that the language does not affect the expressivity of the plain objects. In order to test that condition and prove that the language affects or doesn't affect the expressivity of plain objects we will need first to find two *type-independent languages*.

**Definition 6.1** (Type-independent languages). *Two languages  $L_1$  and  $L_2$  are type-independent if and only if one of the languages contains a type that cannot be represented by means of a linear combination of any other type of the other language.*

For example, let's take Java  $L_1$  and Rust  $L_2$ , examples (a) and (c) from Figure 6.5. Rust contains the type *Either*  $< A, B >$ , this type allows the type  $A$  or  $B$  and when accessed is not an *Either* is either  $A$  or  $B$ . In Java there is no *Either* type, and someone can say that

we could achive a similar type by using inheritance and classes composition. But at the end when accessed the type would be the type of the upper class. **Therefore Java and Rust are type-independent languages.**

Now in order to see if the expressivity depends on the types of a language let's assing values to Java and Rust by using the same *Either*  $< A, B >$  type. As can be see in [Figure 6.7](#) Java does not allow to express the same as Rust is expressing in this example. And therefore we can conclude that:

**Partial Conclusion 6.3.** *The expresivity of plain object is stronggly related to the build-in types that the programming language in which they are coded provides.*

Person Rust Struct	Person Java Object
6 struct Person {	8 // Imports...
7     name: String,	9 public class Person {
8     knows: List<Person>,	10     private String name;
9     owningPet: Either<Dog,Cat>,	11     private List<Person> knows;
10 }	12     private Pet owningPet;
	13     // Constructor...
	14     // Getters and Setters...
	15 }

**Figure 6.7:** Rust struct modeling a **Person** to the left. And the most similar approximation in Java to te right. In the Java approximation the Pet class is an interface that it is inherited by the Cat and Dog classes, that way we allow to store in the variable `owningPet` values of type Cat and Dog.

6.2.3 Plain Objects Expressivity Generalization

6.3 Shape Expressions and Plain Objects Expressivity Comparison

Previous section cover the expressivity of Shape Expressions and Plain Objects, in this section we compare both expressivities and expose if both expresivities are fully compatible or not.

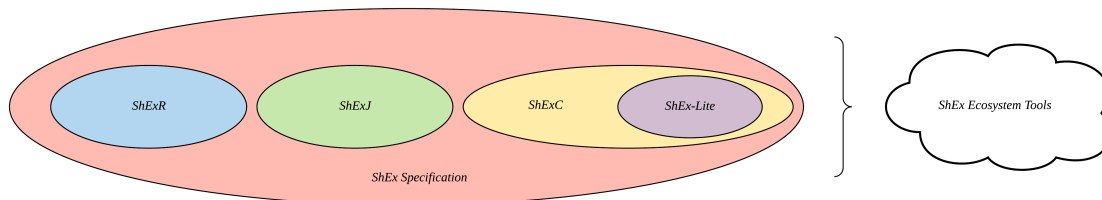
Generalizar  
un PO viene  
ser una lista  
propiedad v  
donde lo im  
tante son lo  
de cada leng  
sin embargo  
tipos princip  
que todos lo  
PO cumplen  
los Strings,   
números ent  
y de coma fl  
referencias a  
objetos y lis  
los tipos ant

Decir los PO  
sentan más  
stricciones q

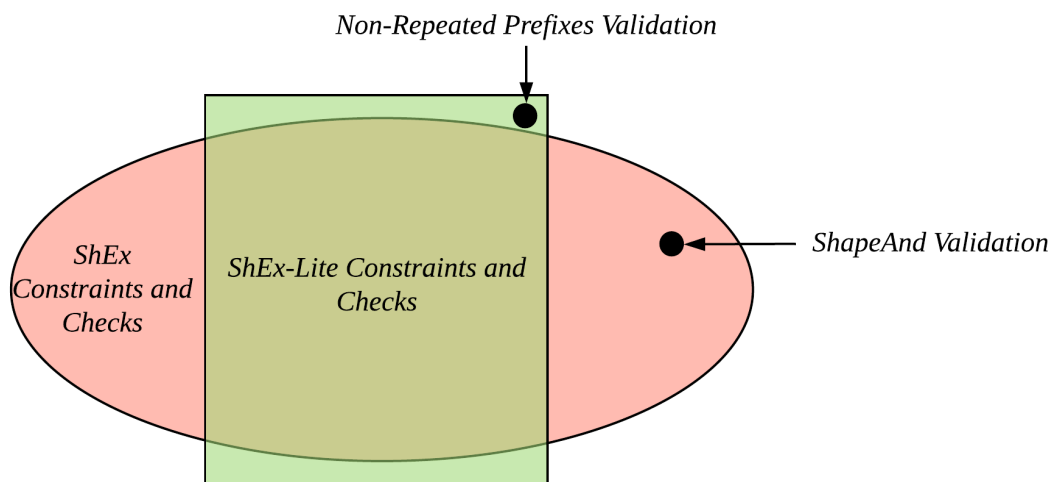
# CHAPTER 7

## Proposed Translator

Para solucionar el problema expuesto en el punto 4, con toda esa información proponemos una sintaxis basada en shexc que representa el subconjunto que podemos representar con PO. Dicha sintaxis a demás al ser un subconjunto está validada semantica y sintácticamente de forma que tiene más restricciones que shex normal, de forma que se aseguran las integridades del lenguaje y se propone un sistema de errores "actualizado". Para esta sintaxis ofrecemos tb un compilador diseñado como una API que compila los esquemas que se realizan en esta sintaxis y los traduce al lenguaje de programación objetivo que se desee.



**Figure 7.1:** Mental model of ShEx-Lite in the existing ShEx syntaxes context. From this model we can see that Shex-Lite is in fact an strictly subset of ShExC, which follows the ShEx Specification. And therefore ShEx-Lite will also follow that expecification, which automatically enables ShEx-Lite schemas to be used in any other existing ShEx tool.



**Figure 7.2:** Constraints and checks context diagram for ShEx-Lite and ShEx.

## 7.1 Structure

## 7.2 Generated Obejcts

Person.shexl	Person.java
<pre> 11 # Prefixes... 12 :Person { 13     :name xsd:string ; 14     :knows @:Person * 15 }</pre>	<pre> 16 // Imports... 17 public class Person { 18     private String name; 19     private List&lt;Person&gt; knows; 20     // Constructor... 21     // Getters and Setters... 22 }</pre>

**Figure 7.3:** Schema modeling a Person in shexl syntax to the left. And the ShEx-Lite generated code in Java to the right.

## Part III

# Project Synthesis

# CHAPTER 8

## Evaluation of Results

### 8.1 Methodology

In order to evaluate the proposed solutions to the two questions posed in [Chapter 1](#), the following methodologies has been used.

1. To evaluate how error detection has been improved, compare the number of actual errors found in a form by the tools detected and against the proposed solution.
2. To evaluate how the error information system has been improved, the number of elements that make up the error messages of each existing tool and of our solution is compared. In addition, a survey is carried out on different users familiar with the existing tools.
3. To evaluate to what extent we can translate shapes to domain object models we collect all the existing shapes in github, reduce the set to those that fit the micro compact syntax and try to generate objects for those that are syntactically and semantically valid. In this way we can approximate what percentage we can translate.

# CHAPTER 9

## Planning and Budget

### 9.1 Planning

### 9.2 Budget

# CHAPTER 10

## Conclusions

### 10.1 Future Work



## Part IV

# Annexes and References

# APPENDIX A

## ShEx Micro Language

### A.1 Syntax Specification

### A.2 Lexical Specification

# APPENDIX B

## ShEx-Lite Antlr Grammar

### B.1 Syntax Specification

### B.2 Lexical Specification

# APPENDIX C

## Project Communications

### C.1 Open Source Community

### C.2 Scientific Disclosure

### C.3 Dublin Core Meetings

# Bibliography

- [1] dotnet/roslyn. Library Catalog: [github.com](https://github.com).
- [2] Domain-specific language, March 2020. Page Version ID: 945660040.
- [3] Entorno de desarrollo integrado, May 2020. Page Version ID: 126109294.
- [4] Lexical analysis, May 2020. Page Version ID: 955735363.
- [5] List of object-oriented programming languages, April 2020. Page Version ID: 950660258.
- [6] Optimizing compiler, April 2020. Page Version ID: 949660274.
- [7] Programming language, May 2020. Page Version ID: 958722580.
- [8] Turing completeness, May 2020. Page Version ID: 954916159.
- [9] Tim Berners-Lee et al. Semantic web road map, 1998.
- [10] Tim Berners-Lee, James Hendler, and Ora Lassila. The semantic web. *Scientific american*, 284(5):34–43, 2001.
- [11] Iovka Boneva, Jérémie Dusart, Daniel Fernández Alvarez, and Jose Emilio Labra Gayo. Semi automatic construction of shex and shacl schemas. 2019.
- [12] Min Chen, Shiwen Mao, and Yunhao Liu. Big data: A survey. *Mobile networks and applications*, 19(2):171–209, 2014.
- [13] Daniel Fernández-Alvarez, Jose Emilio Labra-Gayo, and Herminio García-González. Inference and serialization of latent graph schemata using shex, 2016.
- [14] Robert W Floyd. Syntactic analysis and operator precedence. *Journal of the ACM (JACM)*, 10(3):316–333, 1963.
- [15] Martin Fowler. *Analysis patterns: reusable object models*. Addison-Wesley Professional, 1997.
- [16] Herminio Garcia-Gonzalez, Daniel Fernandez-Alvarez, and Jose Emilio. Shexml: An heterogeneous data mapping language based on shex.
- [17] Tom Heath and Christian Bizer. Linked data: Evolving the web into a global data space. *Synthesis lectures on the semantic web: theory and technology*, 1(1):1–136, 2011.

- [18] Jose Emilio Labra-Gayo, Herminio García-González, Daniel Fernández-Alvarez, and Eric Prud'hommeaux. Challenges in rdf validation. In *Current Trends in Semantic Web Technologies: Theory and Practice*, pages 121–151. Springer, 2019.
- [19] Jose Emilio Labra Gayo, Eric Prud'Hommeaux, Iovka Boneva, and Dimitris Kontokostas. Validating rdf data. *Synthesis Lectures on Semantic Web: Theory and Technology*, 7(1):1–328, 2017.
- [20] Mark Levene and George Loizou. A graph-based data model and its ramifications. *IEEE Transactions on Knowledge and Data Engineering*, 7(5):809–823, 1995.
- [21] Frank Manola, Eric Miller, Brian McBride, et al. Rdf primer. *W3C recommendation*, 10(1-107):6, 2004.
- [22] Terence J. Parr and Russell W. Quong. Antlr: A predicated-ll (k) parser generator. *Software: Practice and Experience*, 25(7):789–810, 1995.
- [23] Eric Prud'hommeaux, Jose Emilio Labra Gayo, and Harold Solbrig. Shape expressions: an rdf validation and transformation language. In *Proceedings of the 10th International Conference on Semantic Systems*, pages 32–40, 2014.
- [24] Arthur G Ryman, Arnaud Le Hors, and Steve Speicher. Oslc resource shape: A language for defining constraints on linked data. *LDOW*, 996, 2013.
- [25] Eric Van der Vlist. *Relax ng: A simpler schema language for xml*. " O'Reilly Media, Inc.", 2003.
- [26] Peter Wegner. Concepts and paradigms of object-oriented programming. *ACM Sigplan Ops Messenger*, 1(1):7–87, 1990.