

Formulating RDF Constraints and Validating RDF Data Generically

Thomas Bosch¹, Erman Acar², and Kai Eckert²

¹ GESIS – Leibniz Institute for the Social Sciences, Germany
thomas.bosch@gesis.org,

² University of Mannheim, Germany
{erman,kai}@informatik.uni-mannheim.de

Abstract. More and more domains adopt Linked Data principles and seek for ways to ensure data quality by validating RDF data according to predefined constraints, particularly as many communities are used to have this functionality in the XML world. XML Schema is the standard constraint language to formulate constraints on and to validate XML documents. For RDF, however, there are several candidates for languages to express constraints on RDF data - but no standard one meeting the majority of RDF validation use cases. We identified requirements to formulate constraints and to validate RDF data and mapped each requirement to an RDF constraint type. The majority of these constraint types can be represented in description logics providing their logical underpinning. The **contributions** of this paper are: (1) we provide a basic terminology and classification system for RDF constraints, (2) we developed a vocabulary to describe constraints of any constraint type generically, (3) we show how to transform constraints, expressed by any constraint language α , into generically expressed constraints and into constraints represented by any other constraint language β and discuss why these transformations are useful, and (4) we explain how to overcome the necessity to implement the validation of each constraint type for multiple constraint languages.

Keywords: RDF Validation, RDF Constraints, RDF Validation Requirements, Linked Data, Semantic Web

1 Introduction

In 2013, the W3C organized the RDF Validation Workshop³, where experts from industry, government, and academia discussed first use cases for RDF constraint⁴ formulation and RDF data validation. In 2014, two working groups on

³ <http://www.w3.org/2012/12/rdf-val/>

⁴ For simplicity reasons, we use the terms *constraint types/constraints* instead of *RDF constraint types/RDF constraints* in the rest of the paper

Thomas: feedback: the major problem is semantics. OWL basically work in different semantics than others. And their semantics are formal. Closed world vs. Open World is also a part of this.

RDF validation have been established to develop a language to express constraints on RDF data: the W3C RDF Data Shapes working group⁵ and the DCMI RDF Application Profiles task group⁶. Bosch and Eckert [1] collected the findings of these working groups and initiated a database of RDF validation requirements which is available for contribution at <http://purl.org/net/rdf-validation>. The intention is to collaboratively collect case studies, use cases, requirements, and solutions regarding RDF validation in a comprehensive and structured way. Bosch et al. identified 74 requirements to formulate constraints; each of them corresponding to a constraint type which may be expressed by at least one existing constraint language [3].

Recently, RDF validation as a research field gained speed due to common needs of data practitioners. A typical example is the library domain that co-developed and adopted Linked Data principles very early. For libraries, the common description of resources are key business and they have a long tradition in developing and using interoperable data formats. While they embrace the openness of Linked Data and the data modeling principles provided by RDF, the data is still mostly represented in XML and this is unlikely to change soon. Among the reasons for the success of XML is the possibility to formulate fine-grained constraints to be met by the data and to validate the data according to these constraints using powerful systems like DTDs, XML Schemas, RELAX NG, or Schematron. A typical example is the definition of a library record describing a book. There are clear rules which information has to be available to describe a book properly (required fields, like a title), but also how information like an ISBN number is properly represented. Libraries seek to make their own data reusable for general purposes, but also to enrich and interlink their own data. Checking if third-party data meets own requirements or validating existing data according to new needs for a Linked Data application are among common use cases for RDF validation.

2 Motivation

There is not the one language meeting the majority of requirements to formulate constraints. As there are more than 10 candidate languages with different syntaxes and semantics which can be used to express constraints, the language selection heavily depends on the individual use case. The five most promising constraint languages on being the standard are *Description Set Profiles (DSP)*, *Resource Shapes (ReSh)*, *Shape Expressions (ShEx)*, the *SPARQL Inferencing Notation (SPIN)*, and the *Web Ontology Language (OWL 2)*. The idea behind this paper is to achieve complete RDF validation requirements coverage by representing any constraint type (corresponding to a requirement to formulate constraints) in a generic way using a lightweight vocabulary consisting of only a few terms (see section 3).

⁵ <http://www.w3.org/2014/rds/charter>

⁶ <http://wiki.dublincore.org/index.php/RDF-Application-Profiles>

Cardinality restrictions, e.g., can be expressed either generically (*generic constraints*) by *description logics* (*DL*) or specifically (*specific constraints*) by domain-specific constraint languages such as *OWL 2*, *DSP*, *ShEx*, *ReSh*, and *SPIN*. The knowledge representation formalism *DL*, with its well-studied theoretical properties, provides the foundational basis to express constraints in a generic way. For that reason, we map constraint types to *DL* statements. With *minimum qualified cardinality restrictions* (*R-75*, *R-81*)⁷, researchers from the library domain may restrict that *publications* must have at least one *author* which must be a *person*. This constraint can be represented generically (*generic constraint*) using *DL* ($\text{Publication} \sqsubseteq \geq 1 \text{ author.Person}$), but also specifically (*specific constraint*) by multiple constraint languages:

```

1 # OWL 2:
2 Publication
3     a owl:Restriction ;
4     owl:minQualifiedCardinality 1 ;
5     owl:onProperty author ;
6     owl:onClass Person .
7
8 # ShEx:
9 Publication { author @Person{1, } }
10
11 # ReSh:
12 Publication a rs:ResourceShape ; rs:property [
13     rs:name "author" ; rs:propertyDefinition author ;
14     rs:valueShape Person ;
15     rs:occurs rs:One-or-many ; ] .
16
17 # DSP:
18 [
19     a dsp:DescriptionTemplate ;
20     dsp:resourceClass Publication ;
21     dsp:statementTemplate [ a dsp:NonLiteralStatementTemplate ;
22         dsp:minOccur 1 ; dsp:maxOccur "infinity" ;
23         dsp:property author ;
24         dsp:nonLiteralConstraint [ a dsp:NonLiteralConstraint ;
25             dsp:valueClass Person ] ] ] .

```

As each *publication* must have at least one *author* which must be a *person* and as the book *The-Lord-Of-The-Rings* is a *publication* ($\text{rdf:type}(\text{The-Lord-Of-The-Rings}, \text{Publication})$), *The-Lord-Of-The-Rings* must have at least one *author* relationship to a *person*. Constraint violations are raised either if *The-Lord-Of-The-Rings* does not have any *author* relationship, or if *The-Lord-Of-The-Rings* has an *author* which is not a *person* ($\text{author}(\text{The-Lord-Of-The-Rings}, \text{Tolkien})$, $\text{rdf:type}(\text{Tolkien}, \text{Hobbit})$), or if *The-Lord-Of-The-Rings* has an *author* for which no class is assigned ($\text{author}(\text{The-Lord-Of-The-Rings}, \text{Tolkien})$). The data is valid, however, if *The-Lord-Of-The-Rings* is connected to a *person* via the property *author* ($\text{author}(\text{The-Lord-Of-The-Rings}, \text{Tolkien})$, $\text{rdf:type}(\text{Tolkien}, \text{Person})$).

As there is no standard way to formulate constraints, semantically equivalent *cardinality restrictions* may be represented by multiple constraint languages with different syntax and semantics. This causes confusion and weakens the common

⁷ Requirements/constraint types are uniquely identified by alphanumeric technical identifiers like *R-75-MINIMUM-QUALIFIED-CARDINALITY-ON-PROPERTIES*

Thomas: Erman thinks formal generic semantics should be defined

understanding between several parties about the semantics of particular constraints and therefore how to ensure high data quality. Thus, when choosing constraint language α to express a constraint of an arbitrary constraint type, it should be possible to transform this constraint into a semantically equivalent constraint formulated by any other constraint language β (see section 4). This is important in order to enhance the interoperability of constraint languages and to resolve ambiguities in the communication of RDF data producers and consumers, as constraint transformations avoid the necessity to understand several constraint languages.

We provide a validation environment (available at <http://purl.org/net/rdfval-demo>) which enables to validate RDF data according to constraints expressed by diverse constraint languages. The next step, however, is to extend the **RDF Validator** by validating constraints of any constraint type represented by any constraint language and to validate semantically equivalent *specific constraints*, expressed by different constraint languages, in exactly the same way, i.e., validation is independent from the used constraint language (see section 5). We are able to offer an already implemented validation mechanism for any *specific constraint* out-of-the-box, (1) if we implement the validation of the corresponding *generic constraint type* once and (2) if we transform semantically equivalent *specific constraints* into generically expressed constraints.

This paper aims to address two main **audiences**: (1) RDF data providers and consumers seeking how to ensure high quality (meta)data and thus for ways to enhance the interoperability and common understanding of constraints expressible by multiple constraint languages and (2) RDF practitioners thinking of how to provide aligned implementations for the validation of both existing and newly developed constraint languages. The main **contributions** of this paper are: (1) we provide a basic terminology and classification system for RDF constraints (section 3), (2) we developed a vocabulary to describe constraints of any RDF constraint type generically (section 3), (3) we show how to transform constraints, expressed by any constraint language α , into generically expressed constraints and into constraints represented by any other constraint language β and discuss why these transformations are useful (section 4), and (4) we explain how to overcome the necessity to implement the validation of each constraint type for multiple constraint languages (section 5).

3 A Vocabulary to Describe RDF Constraints Generically

In order to develop a vocabulary to describe constraints of any constraint type generically, it is needed to define the basic terminology for the formulation of constraints and to classify them first. A **constraint language** is a language which is used to formulate constraints. The W3C Data Shapes working group defines **RDF constraint** as a component of a schema what needs to be satisfied⁸. We identified four dimensions to classify constraints:

⁸ <https://www.w3.org/2014/data-shapes/wiki/Glossary>

- **Universality:** *specific constraints* vs. *generic constraints*
- **Complexity:** *simple constraints* vs. *complex constraints*
- **Context:** *property constraints* vs. *class constraints*
- **DL Expressivity:** *constraints expressible in DL* vs. *constraints not expressible in DL*

As there are already five promising constraint languages, our purpose is not to invent a new constraint language. We rather developed a very simple lightweight vocabulary (only three classes, three object properties, and three data properties) which is universal enough to describe constraints of any constraint type expressible by any constraint language (see the conceptual model in figure 1). We call this vocabulary the **RDF Constraints Vocabulary (RDF-CV)**⁹.

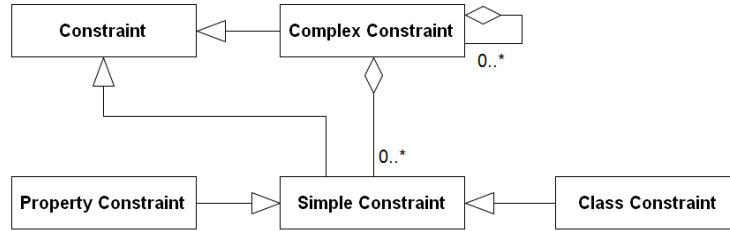


Fig. 1. *RDF Constraints Vocabulary (RDF-CV)* - conceptual model

Constraints can be expressed either generically (**generic constraints**) by the *RDF-CV* or specifically (**specific constraints**) by domain-specific constraint languages. As the *RDF-CV* describes constraints generically, the vocabulary does not distinguish constraints according to the dimension *universality*. The majority of the constraint types can be expressed in *DL*. In contrast, there are constraint types which cannot be expressed in *DL*, but are also representable by the *RDF-CV* (see section 6). **Complex constraints** encompass **simple constraints** (atomic constraints) and/or further *complex constraints*, i.e. *DL* statements representing *complex constraints* are created out of *DL* statements standing for atomic and composed constraints (if expressible in *DL*). Simple constraints may be applied to either properties (**property constraints**) or classes (**class constraints**). The *RDF-CV* does not contain any terms standing for *simple* and *complex constraints*, since *context classes* of *simple constraints* are just reused within *complex constraints* (*simple constraints* associated with *context classes* must hold for individuals of these *context classes*). As a consequence, the distinction of *property* and *class constraints* is sufficient to describe constraints of all possible constraint types.

3.1 Simple Constraints (Expressible in *DL*)

Sub-classes of **simple constraints** are **property constraints** and **class constraints** (see the *RDF-CV* implementation model in figure 2). For both *property*

⁹ Available at: <https://github.com/boschthomas/RDF-Constraints-Vocabulary>

and *class constraints* a *context class*, a list of *classes*, the *constraining element*, and the *constraining value* can be stated. Lists of *left* and *right properties* can only be specified for *property constraints*.

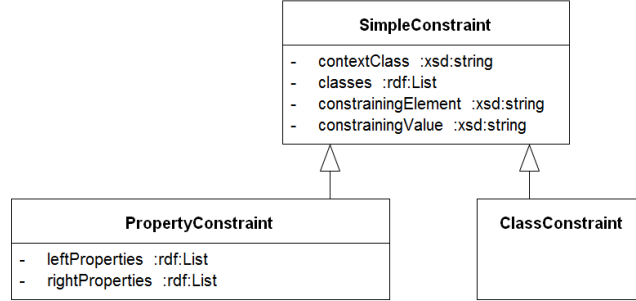


Fig. 2. *RDF Constraints Vocabulary (RDF-CV)* - implementation model

A *simple constraint* holds for all individuals of a ***context class***. Consider, e.g., the following *minimum qualified cardinality restriction* on the object property *author*, which restricts *publications* to have at least one *author* which is a *person*:

$$Publication \sqsubseteq_{\geq 1} author.Person$$

The cardinality restriction is a constraint on the property *author* and is therefore classified as a *property constraint* which holds for all individuals of the class *Publication* (see table 1 for the mapping to *RDF-CV*):

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	Publication	author	-	Person	\geq	1

Table 1. Property Constraint - Cardinality Restriction

The ***constraining element*** indicates the actual type of constraint like *DL* concept and role constructors, (in)equality, and further keywords for constraint types which cannot be expressed in *DL* (e.g. regular expressions and constraining facets). In some cases, a constraint is only complete when in addition to the *constraining element* a ***constraining value*** is stated. The cardinality restriction $\geq 1 \text{ author} . Person$, e.g., constructs an anonymous class of all individuals having at least one *author* relationship to *persons*. The *constraining element* of this *property constraint* is the *DL* at-least restriction \geq and the *constraining value* is 1. *Simple constraints* may refer to a list of ***classes***. The qualified cardinality restriction above, e.g., refers to the class *Person*, i.e. restricts the objects of the property *author* to be of the class *Person*.

For *property constraints*, ***left*** and ***right property lists*** are specified. The assignment of properties to these lists happens relative to the *constraining element* - which may be an operator (e.g., \sqsubseteq in case of *object property paths*). *Object Prop-*

erty Paths¹⁰ (also called *complex role inclusion axioms* in DL) express that, if an individual x is connected by a sequence of object properties with an individual y , then x is also related to y by a particular object property. As *Stephen-Hawking* is the *author* of the book *A-Brief-History-Of-Time* (`authorOf(Stephen-Hawking, A-Brief-History-Of-Time)`) whose *genre* is *Popular-Science* (`genre(A-Brief-History-Of-Time, Popular-Science)`), the following *object property path* infers that *Stephen-Hawking* is an *author* of the *genre Popular-Science* (`authorOfGenre(Stephen-Hawking, Popular-Science)`):

$$\text{authorOf} \circ \text{genre} \sqsubseteq \text{authorOfGenre}$$

Thus, when mapped to the *RDF-CV* (see table 2), the properties *authorOf* and *genre* are on the left side of the *constraining element* \sqsubseteq , and the property *authorOfGenre* is on the right side:

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	T	authorOf, genre	authorOfGenre	T	\sqsubseteq	-

Table 2. Property Constraint - Object Property Paths

3.2 Simple Constraints (Not Expressible in DL)

There are *simple constraints* which cannot be expressed in DL such as *literal pattern matching*, *literal value comparison*, *literal ranges*, *default values*, and *language tag cardinality* [3].

*Literal pattern matching*¹¹ restricts literals to match given patterns. The universal restriction `Book \sqsubseteq \forall identifier.ISBN`, e.g., ensures that *books* can only have valid ISBN identifiers, i.e. strings that match a given regular expression. The restriction of the datatype *ISBN*, however, cannot be expressed in DL, but OWL 2 DL can be used anyway to express the *literal pattern matching* constraint:

```

1 ISBN a RDFS:Datatype ; owl:equivalentClass [ a RDFS:Datatype ;
2   owl:onDatatype xsd:string ; owl:withRestrictions ([ xsd:pattern "~\d{9}[\d|X]\$" ]) ] .

```

The first OWL 2 axiom explicitly declares *ISBN* to be a datatype. The second OWL 2 axiom defines *ISBN* as an abbreviation for a datatype restriction on *xsd:string*. The datatype *ISBN* can be used just like any other datatype like in the universal restriction above. The *literal pattern matching* constraint validates *ISBN* literals according to the stated regular expression causing constraint violations for triples not matching the regular expression. In table 3, the *simple constraint* (*literal pattern matching*) and the *complex constraint* (*universal restriction*) are mapped to the *RDF-CV*.

Thomas: TODO: delete the slash in front of the dollar sign from the output of the code fragment above (in the stated pattern)?

¹⁰ Corresponds to the requirement *R-55-OBJECT-PROPERTY-PATHS*

¹¹ Corresponds to *R-44-PATTERN-MATCHING-ON-RDF-LITERALS*

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	Book	identifier	-	ISBN	\forall	-
class	ISBN	-	-	xsd:string	regex	'^\d{9}[\d X]\$',

Table 3. Simple Constraint (Not Expressible in DL)

The *literal pattern matching* constraint introduces the new *constraining element regex*. For this new *constraining element*, the validation must be implemented. This has to be done once for each *generic constraint type* which cannot be expressed in DL.

3.3 Complex Constraints

Complex constraints encompass *simple constraints* and/or *complex constraints*. *Exclusive or* is a logical operation that outputs true whenever both inputs differ (one is true, the other is false). The *Context-Specific Exclusive OR of Property Groups*¹² is a *complex constraint* restricting individuals of a given context (here a class) to contain only one of multiple property groups. Publications, e.g., are either identified by an ISBN and a title (for books) or by an ISSN and a title (for periodical publications), but it should not be possible to assign both identifiers to a given publication, which can be represented in ShEx:

```

1 Publication {
2   ( isbn string , title string ) |
3   ( issn string , title string ) }
```

As, e.g., the *The-Great-Gatsby* is a *publication* (`rdf:type(The-Great-Gatsby, Publication)`) having an *ISBN* number (`isbn(The-Great-Gatsby, '978-0241965672')`) and a *title* (`title(The-Great-Gatsby, 'The Great Gatsby')`) and not an additional *ISSN* identifier, the data is valid. The *complex constraint* above is mapped to DL as follows:

$$\begin{aligned}
\text{Publication} &\sqsubseteq (\neg E \sqcap F) \sqcup (E \sqcap \neg F) \\
E &\equiv A \sqcap B \\
F &\equiv C \sqcap D \\
A &\sqsubseteq \geq 1 \text{ isbn.string} \sqcap \leq 1 \text{ isbn.string} \\
B &\sqsubseteq \geq 1 \text{ title.string} \sqcap \leq 1 \text{ title.string} \\
C &\sqsubseteq \geq 1 \text{ issn.string} \sqcap \leq 1 \text{ issn.string} \\
D &\sqsubseteq \geq 1 \text{ title.string} \sqcap \leq 1 \text{ title.string}
\end{aligned}$$

The *complex constraint* (expressed in DL) is mapped to a *generic constraint* using the *RDF-CV* (see table 4).

Even though, tools may generate *generic constraints* automatically, this *complex constraint* is composed of many other *complex* (e.g. minimum and maximum qualified cardinality restrictions) and *simple constraints* (e.g. constraints on sets). As *exact (un)qualified cardinality restrictions (=n)* and *exclusive or (XOR)* are frequently used *complex constraints*, we propose using them as *simple*

¹² Corresponds to R-13-DISJOINT-GROUP-OF-PROPERTIES-CLASS-SPECIFIC

c. type	context class	left p. list	right p. list	classes	c. element	c. value
class	Publication	-	-	$\neg E \sqcap F, E \sqcap \neg F$	\sqcup	-
class	$\neg E \sqcap F$	-	-	$\neg E, F$	\sqcap	-
class	$E \sqcap \neg F$	-	-	$E, \neg F$	\sqcap	-
class	$\neg E$	-	-	E	\sqcup	-
class	E	-	-	A, B	\sqcup	-
class	$\neg F$	-	-	F	\sqcup	-
class	F	-	-	C, D	\sqcup	-
class	A	-	-	$A1, A2$	\sqcup	-
property	A1	isbn	-	string	\sqcup	1
property	A2	isbn	-	string	\sqcup	1
class	B	-	-	$B1, B2$	\sqcup	-
property	B1	title	-	string	\sqcup	1
property	B2	title	-	string	\sqcup	1
class	C	-	-	$C1, C2$	\sqcup	-
property	C1	issn	-	string	\sqcup	1
property	C2	issn	-	string	\sqcup	1
class	D	-	-	$D1, D2$	\sqcup	-
property	D1	title	-	string	\sqcup	1
property	D2	title	-	string	\sqcup	1

Table 4. Complex Constraints

constraints - in terms of syntactic sugar. As a consequence, the *Context-Specific Exclusive OR of Property Groups* constraint above is represented as a *generic constraint* more intuitively and concisely (see table 5).

c. type	context class	left p. list	right p. list	classes	c. element	c. value
class	Publication	-	-	E, F	XOR	-
class	E	-	-	A, B	\sqcap	-
class	F	-	-	C, D	\sqcap	-
property	A	isbn	-	string	=	1
property	B	title	-	string	=	1
property	C	issn	-	string	=	1
property	D	title	-	string	=	1

Table 5. Simplified Complex Constraints

3.4 Complex Constraints as Simple Constraints (Syntactic Sugar)

Almost 15 percent of all RDF constraint types are *complex constraints* which can be simplified and therefore formulated as *simple constraints* when using them in terms of syntactic sugar (see section 6).

There are three forms of *OWL RBox axioms*: *role inclusions*, *equivalence* and *disjointness*. OWL provides a variety of others, namely *role transitivity*, *symmetry*, *asymmetry*, *reflexivity* and *irreflexivity*. These are sometimes considered as basic axiom types in DL as well, using some suggestive notation such as **Trans(ancestorOf)** to express that the role *ancestorOf* is transitive. However, such axioms are just syntactic sugar; all role characteristics can be expressed using the basic features of DL. The *Irreflexive Object Properties* constraint type¹³, e.g., restricts that an object property is irreflexive if it is never locally reflexive, i.e. no individual is connected by the object property to itself [4]. The irreflexivity constraint $\top \sqsubseteq \neg \exists \text{authorOf.Self}$ states that individuals cannot be authors of

Thomas: transitive property example from the library domain

¹³ Corresponds to *R-60-IRREFLEXIVE-OBJECT-PROPERTIES*

themselves. When mapped to the *RDF-CV* (see table 6), the *complex constraint* aggregates three *simple constraints* (one *property* and two *class constraints*).

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	\exists authorOf.Self	marriedTo	-	Self	\exists	-
class	$\neg \exists$ authorOf.Self	-	-	\exists authorOf.Self	\neg	-
class	\top	-	-	$\top, \neg \exists$ authorOf.Self	\sqsubseteq	-

Table 6. Irreflexive Object Properties as Complex Constraints

When using the OWL RBox axiom *role irreflexivity* as syntactic sugar, the *complex constraint* transforms into a *simple property constraint* with exactly the same semantics (see table 7).

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	\top	authorOf	-	-	irreflexive	-

Table 7. Irreflexive Object Properties as Simple Constraints

The *Primary Key Properties*¹⁴ constraint type is often useful to declare a given (datatype) property as the "primary key" of a class, so that a system can enforce uniqueness. Books, e.g., are uniquely identified by their ISBN, i.e. the property *isbn* is inverse functional (funct *isbn*⁻), which can be mapped to the *RDF-CV* (see table 8).

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	\top	isbn ⁻	isbn ⁻	-	inverse	-
property	Book	isbn ⁻	-	-	functional	-

Table 8. Primary Key Properties as Complex Constraints

Keys, however, are even more general, i.e. a generalization of inverse functional properties [6]. A key can be a datatype property, an object property, or a chain of properties. For these generalization purposes, as there are different sorts of keys, and as keys can lead to undecidability, DL is extended with *key boxes* and a special construct *keyfor* [5]. When using the *keyfor* DL construct (*isbn keyfor Book*), the *complex constraint* can be expressed by only one *simple property constraint* (see table 9 for the mapping to the *RDF-CV*).

4 How to Improve the Interoperability of Constraints

RDF data providers have to ensure that their data conforms to constraints which are part of explicitly or implicitly concluded contracts between data providers and consumers. As there is no standard way, these constraints may be represented by a variety of possible constraint languages (there are even more than the five promising ones¹⁵). There is an obvious need to transform any *specific constraint* (sc_α) (expressed by any constraint language α) into any other *specific constraint* (sc_β) (expressed by any constraint language β) - when they have the

¹⁴ Corresponds to *R-226-PRIMARY-KEY-PROPERTIES*

¹⁵ i.a. Bibframe, DQTP, Pellet ICV, RDFUnit, SPARQL, Stardog ICV

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	Book	isbn	-	-	keyfor	-

Table 9. Primary Key Properties as Simple Constraints

same semantics. We use the *RDF-CV* as an intermediate representation vocabulary for these transformations. If we define mappings $(m(sc_\alpha), m(sc_\beta))$ between equivalent *specific constraints* and the corresponding *generic constraint* (gc), we can convert them automatically:

$$gc = m(sc_\alpha) = m(sc_\beta)$$

We do not need to define mappings for each pair of semantically equivalent *specific constraints*. Let's assume we are able to express a constraint type by 5/6/7/8/9/10 constraint languages. Without the intermediate *generic constraint*, we would have to define for each constraint type l mapping for each pair of *specific constraints* (expressed by different constraint languages) - that are $\binom{n}{2}$ mappings (10/15/21/28/36/45 mappings). With an intermediate *generic constraint*, however, we only need to define for each constraint type n mappings (one mapping for each *specific constraint*) from n *specific constraints* to the corresponding *generic constraint* (5/6/7/8/9/10 mappings).

The *existential quantification*¹⁶ `Publication \equiv \exists publisher.Publisher` restricts publications to have at least one publisher and can be expressed by at least 10 constraint languages¹⁷:

```

1  # OWL 2:
2  [ a owl:Restriction ;
3    owl:onProperty publisher ;
4    owl:someValuesFrom Publisher ;
5    rdfs:subClassOf Publication ] .
6
7  # SPARQL:
8  ASK { ?x a Publication . FILTER NOT EXISTS { ?x publisher [ a Publisher ] } . }
9
10 # ShEx:
11 Publication {
12   publisher @Publisher{1,} }
13
14 # DSP:
15 descriptionTemplate a dsp:DescriptionTemplate ;
16   dsp:resourceClass Publication ;
17   dsp:statementTemplate [ a dsp:NonLiteralStatementTemplate ;
18     dsp:minOccur 1 ; dsp:maxOccur "infinity" ;
19     dsp:property publisher ;
20     dsp:nonLiteralConstraint [ a dsp:NonLiteralConstraint ;
21       dsp:valueClass Publisher ] ] .
22
23 # ReSh:
24 Publication a rs:ResourceShape ; rs:property [
25   rs:propertyDefinition :publisher ;
26   rs:valueShape :Publisher ;
27   se:min 1 ; se:maxundefined true ; ] .
28

```

¹⁶ Corresponds to *R-86-EXISTENTIAL-QUANTIFICATION-ON-PROPERTIES*

¹⁷ at least by OWL2, SPARQL, ShEx, Bibframe, DSP, ReSh, SPIN, DQTP, Pellet IVC, Stardog IVC

```

29 # SPIN:
30 CONSTRUCT {
31   _:cv a spin:ConstraintViolation ; spin:violationRoot ?subject ;
32   rdfs:label ?violationMessage ; spin:violationPath ?property . } WHERE
33 { ?subject a ?classSubject . FILTER NOT EXISTS { ?subject ?property [ a ?classObject ] } . }

```

First, the DL statement $\text{Publication} \equiv \exists \text{ publisher.Publisher}$ is mapped to a *simple generic constraint* according to the *RDF-CV* (see table 10).

c. type	context class	left p. list	right p. list	classes	c. element	c. value
property	Publication	publisher	-	Publisher	\exists	

Table 10. Existential Quantification as Simple Constraints

The next step is to define one mapping for each *specific constraint* (sc_{OWL2} , sc_{ReSh} , sc_{ShEx} , ...) to the *generic constraint*.

5 How to Provide Constraint Validation Out-Of-The-Box

SPARQL is generally seen as the method of choice to validate RDF data according to certain constraints, although it is not ideal for their formulation. In contrast, OWL 2, DSP, ReSh, and ShEx constraints are comparatively easy to understand. Bosch and Eckert [2] use SPIN as basis to define a validation environment in which the validation of any constraint language¹⁸ can be implemented by representing them in SPARQL (the implementation can be tested at <http://purl.org/net/rdval-demo>). The SPIN engine checks for each resource if it satisfies all constraints (associated with its assigned classes) and generates a result RDF graph containing information about all constraint violations. We are able to offer the validation of all *specific constraints* (expressed by any constraint language) by defining mappings from each *specific constraint* to SPIN, what we have done for all OWL 2¹⁹ and DSP²⁰ constructs and for some constructs of other constraint languages. We do not have to map each *specific constraint* to SPIN and we get identical validation results for semantically equivalent *specific constraints*, (1) if we map semantically equivalent *specific constraints* to the corresponding *generic constraint* and (2) if we map the *generic constraint* to SPIN. As a consequence, we only have to define one SPIN mapping for each constraint type corresponding to a requirement within the RDF validation requirements database. 19 of 74 RDF constraint types can be expressed by at least four constraint languages [3]. This means that there must be 76 (four for each constraint type) implementations to validate 76 *specific constraints* - instead of only 19 implementations (one for each *generic constraint*).

Within our validation environment, we define one SPIN construct template for each *generic constraint*²¹. A SPIN construct template contains a SPARQL

¹⁸ The only limitation is that constraint languages must be represented in RDF

¹⁹ OWL to SPIN: <https://github.com/boschthomas/OWL2-SPIN-Mapping>

²⁰ DSP to SPIN: <https://github.com/dcmi/DSP-SPIN-Mapping>

²¹ For details about the validation environment see [2]

CONSTRUCT query which generates constraint violation triples (*spin:ConstraintViolation*) indicating the subject (*spin:violationRoot*) and the properties (*spin:violationPath*), causing constraint violations, and the reason why constraint violations have been raised (*rdfs:label*). A SPIN construct template creates constraint violation triples if all triple patterns within the SPARQL WHERE clause match. The SPARQL CONSTRUCT query for validating *universal quantifications*²², e.g. *authors of publications must be persons* ($\text{Publication} \equiv \forall \text{author}.\text{Person}$), is shown below:

```

1 CONSTRUCT {
2   [ a spin:ConstraintViolation ; spin:violationRoot ?subject ;
3     rdfs:label ?violationMessage ; spin:violationPath ?lp1 . ] }
4 WHERE {
5   [ a gclo:PropertyConstraint ;
6     gclo:contextClass ?cc ;
7     gclo:leftProperties ( ?lp1 ) ;
8     gclo:classes ( ?c1 ) ;
9     gclo:constrainingElement "universal quantification" ] .
10  ?subject a ?cc .
11  ?subject ?lp1 ?o .
12  FILTER NOT EXISTS { ?o a ?c1 } . }

```

We provide a mapping from the *RDF-CV* to SPIN²³ in order to automatically validate each type of *generic constraint* (*simple*, *complex*, *property*, and *class constraints*). Part of future work will be (1) to provide a GUI which generates *generic constraints* automatically according to inputs of domain experts who are not familiar with the formulation of constraints, (2) to offer bidirectional transformations between *specific constraints* (expressed by multiple constraint languages) and *generic constraints*, and (3) to provide translations between semantically equivalent *specific constraints* (expressed by any constraint language) by using *generic constraints* as an intermediate transformation step.

6 Evaluation

We evaluated to which extend the most promising five constraint languages fulfill each of the overall 74 requirements to formulate RDF constraints; each of them corresponding to a constraint type. We recently published a technical report²⁴ (serving as appendix of this paper) in which we explain each requirement (constraint type) in detail and give examples for each (represented by listed constraint languages). This technical report also contains mappings to DL to logically underpin each requirement and to determine which DL constructs are needed to express each constraint type [3]. If a constraint type can be expressed in DL, we added the mapping to DL and to the *generic constraint*. If a constraint type cannot be expressed in DL, we only added the mapping to the *generic constraint*. Therefore, we show that each constraint type can be mapped to a *generic constraint*. Table 11 shows the absolute numbers and the relative percentages for

²² Corresponds to *R-91-UNIVERSAL-QUANTIFICATION-ON-PROPERTIES*

²³ RDF-CV to SPIN: <https://github.com/boschthomas/RDF-CV-2-SPIN>

²⁴ Available at: <http://arxiv.org/abs/1501.03933>

each of the three dimensions to classify constraint types *context*, *complexity*, and *DL expressivity*:

Constraint Classes	#	%
<i>Property Constraints</i>	48	65 (64.86)
<i>Class Constraints</i>	17	23 (22.96)
<i>Property and Class Constraints</i>	9	12 (12.16)
<i>Simple Constraints</i>	46	62 (62.16)
<i>Simple Constraints</i> (Syntactic Sugar)	10	14 (13.51)
<i>Complex Constraints</i>	18	24 (24.32)
<i>DL Expressible</i>	51	69 (68.92)
<i>DL Not Expressible</i>	23	31 (31.08)
Total	74	100

Table 11. Evaluation

Constraints types can be classified as *property constraints* and *class constraints* (classification dimension *context*). Two thirds of the total amount of constraint types are *property constraints* (e.g. *existential* and *universal quantification*), one fifth are *class constraints* (e.g. *allowed values*²⁵ for literals and objects), and approx. 10% are composed of both *property* and *class constraints* (e.g. *context-specific exclusive or of property groups*). Constraint types may be either atomic (*simple constraints*) or created out of *simple* and/or *complex constraints* (*complex constraints*) (classification dimension *complexity*). Almost two thirds of the constraint types are *simple constraints* (e.g. *object property paths*), a quarter are *complex constraints* (e.g. *ranges of literal values*²⁶), and nearly 15 percent are *complex constraints* which can be formulated as *simple constraints* when using them in terms of syntactic sugar (e.g. *(un)qualified cardinality restrictions*). Constraint types can either be *expressible in DL* or *not* (classification dimension *DL expressivity*). The majority - nearly 70% - of the overall constraint types are expressible in DL. *Literal pattern matching* is an example of a constraint type which cannot be represented in DL.

7 Related Work

8 Conclusion and Future Work

Thomas: ToDo: there is too much space before and after DL and table environments / can you resolve that please?

References

1. Thomas Bosch and Kai Eckert. Requirements on rdf constraint formulation and validation. *Proceedings of the DCMI International Conference on Dublin Core and Metadata Applications (DC 2014)*, 2014.

²⁵ Corresponds to *R-30-ALLOWED-VALUES-FOR-RDF-OBJECTS* and *R-37-ALLOWED-VALUES-FOR-RDF-LITERALS*

²⁶ Corresponds to *R-45-RANGES-OF-RDF-LITERAL-VALUES*

2. Thomas Bosch and Kai Eckert. Towards description set profiles for rdf using sparql as intermediate language. *Proceedings of the DCMi International Conference on Dublin Core and Metadata Applications (DC 2014)*, 2014.
3. Thomas Bosch, Andreas Nolle, Erman Acar, and Kai Eckert. Rdf validation requirements - evaluation and logical underpinning. 2015.
4. Markus Krötzsch, Frantisek Simancik, and Ian Horrocks. A description logic primer. *CoRR*, abs/1201.4089, 2012.
5. Carsten Lutz, Carlos Areces, Ian Horrocks, and Ulrike Sattler. Keys, nominals, and concrete domains. *Journal of Artificial Intelligence Research*, 23(1):667–726, June 2005.
6. Michael Schneider. OWL 2 Web Ontology Language RDF-Based Semantics. W3C recommendation, W3C, October 2009.