

# RDF Constraints Classification Ensuring High Quality of Metadata and Data

Thomas Bosch<sup>1</sup>, Benjamin Zapilko<sup>1</sup>, Joachim Wackerow<sup>1</sup>, and Kai Eckert<sup>2</sup>

<sup>1</sup> GESIS – Leibniz Institute for the Social Sciences, Germany

`{firstname.lastname}@gesis.org`,

<sup>2</sup> University of Mannheim, Germany

`kai@informatik.uni-mannheim.de`

**Abstract.** For research institutes, data libraries, and data archives, RDF data validation according to predefined constraints is a much sought-after feature, particularly as this is taken for granted in the XML world. To ensure high quality and trust, metadata and data must satisfy certain criteria - specified in terms of RDF constraints.

In this paper, we propose a system to classify RDF constraints and RDF constraint types, which in most cases correspond to RDF validation requirements. Constraints are instantiated from constraint types in order to validate both metadata and data represented by any vocabulary. Within the context of a complex and complete real world running example within the community around research data for the *social, behavioural, and economic (SBE) sciences*, we prove the claim that the developed classification system perfectly applies for diverse vocabularies. We show how data in rectangular format and metadata on person-level data sets (i.e., data about individuals, households, businesses), aggregated data sets, thesauri, and statistical classifications are represented in RDF and how therefore reused vocabularies are interrelated. We explain how *SBE* (meta)data is validated against constraints to ensure high quality and trust. We exhaustively evaluated the metadata quality of large real world aggregated, person-level, and thesauri data sets (more than 4.2 billion triples and 15 thousand data sets) by means of constraints of the majority of constraint types.

**Keywords:** RDF Validation, RDF Constraints, DDI-RDF Discovery Vocabulary, RDF Data Cube Vocabulary, Rectangular Data, SKOS, XKOS, Linked Data, Semantic Web

## 1 Introduction

For more than a decade, members of the community around research data for the *social, behavioural, and economic (SBE) sciences* have been developing and using a metadata standard (composed of almost twelve hundred metadata fields) known as the *Data Documentation Initiative (DDI)* [11]. *DDI* is an XML format designed to support the dissemination, management, and reuse of the data collected and archived for research purposes. Increasingly, data professionals, data

archives, data libraries, government statisticians (e.g. data.gov, data.gov.uk), and national statistical institutes are very interested in having their data be discovered and used by providing their metadata (e.g. about unemployment rates or income) on the web in form of RDF. Recently, members of the SBE and Linked Data community developed the *DDI-RDF Discovery Vocabulary (Disco)*<sup>3</sup>, an effort to leverage the mature DDI metadata model for the purposes of exposing DDI metadata as resources within the Web of Linked Data.

For data archives, research institutes, and data libraries, RDF data validation according to predefined constraints is a much sought-after feature, particularly as this is taken for granted in the XML world as DDI-XML documents are validated against diverse XSDs<sup>4</sup>. Several approaches exist to meet this requirement, ranging from using *OWL 2* as a constraint language to *SPIN*<sup>5</sup>, a SPARQL-based way to formulate and check constraints. There are also constraint languages like *Shape Expressions*, *Resource Shapes* or *Description Set Profiles* that more or less explicitly address the *SBE* community. In 2013, the W3C organized the RDF Validation Workshop<sup>6</sup>, where experts from industry, government, and academia discussed first use cases for RDF constraint formulation and RDF data validation. In 2014, two working groups on RDF validation have been established to develop a language to express constraints on RDF data: the W3C RDF Data Shapes working group<sup>7</sup> and the DCMI RDF Application Profiles task group<sup>8</sup>.

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. The requirements are classified to better evaluate existing solutions and each requirement is directly mapped to a constraint type which may be expressed by at least one existing constraint language. Bosch and Eckert [2] use SPIN as basis to define a validation environment (available at <http://purl.org/net/rdfval-demo>) in which the validation of any constraint language<sup>9</sup> can be implemented by representing them in SPARQL. The SPIN engine checks for each resource if it satisfies all constraints, which are associated with its assigned classes, and generates a result RDF graph containing information about all constraint violations.

The main **contribution** of this paper is the development of a system to classify RDF constraints and RDF constraint types, which in most cases correspond to RDF validation requirements<sup>10</sup> (section 4). We propose an extensible metric to measure the continuum of severity levels to indicate how serious the

<sup>3</sup> <http://rdf-vocabulary.ddialliance.org/discovery.html>

<sup>4</sup> <http://www.ddialliance.org/Specification/>

<sup>5</sup> <http://spinRDF.org/>

<sup>6</sup> <http://www.w3.org/2012/12/rdf-val/>

<sup>7</sup> <http://www.w3.org/2014/rds/charter>

<sup>8</sup> <http://wiki.dublincore.org/index.php/RDF-Application-Profiles>

<sup>9</sup> The only limitation is that constraint languages must be represented in RDF

<sup>10</sup> For simplicity reasons, we use the terms *constraint types* and *constraints* instead of *RDF constraint types* and *RDF constraints* in the rest of the paper

violation of given constraints is. Constraints are instantiated from constraint types in order to validate both metadata and data represented by any vocabulary. As constraint types are used to define constraints on (meta)data expressed by any vocabulary, the proposed constraint classification can be applied generically, i.e. vocabulary-independent. Within the context of a complex and complete real world running example from the *SBE* domain, we prove the claim that the developed classification system perfectly applies for diverse vocabularies. We describe why RDF validation is important for the *SBE* community (section 2), how data in rectangular format (expressed by the *PHDD* vocabulary) and metadata on person-level data sets (*Disco*), aggregated data sets (*QB*), thesauri (*SKOS*), and statistical classifications (*XKOS*) are represented in RDF, and how therefore reused vocabularies are interrelated (section 3). We explain how *SBE* (meta)data is validated against constraints, instantiated from constraint types organized in the constraint classification system, to ensure high quality of and trust in (meta)data.

For the extensive evaluation, we implemented 53 constraint types by instantiating 212 constraints on *Disco*, *QB*, *SKOS*, *XKOS*, and *PHDD* data sets (section 5). First, several *SBE* domain experts evaluated the correctness (i.e., the gold standard) of all constraints and therefore the generic applicability of the developed constraint classification system. Second, we exhaustively evaluated the metadata quality of large real world aggregated (*QB*), person-level (*Disco*), and thesauri (*SKOS*) data sets (more than 4.2 billion triples and 15 thousand data sets) by means of constraints of the majority of constraint types (section 6).

## 2 Motivation

The data most often used in research within the *SBE* community is *person-level data*, i.e. data collected about individuals, businesses, and households, in form of responses to studies or taken from administrative registers (such as hospital records, registers of births and deaths). The range of person-level data is very broad - including census, education, and health data as well as all types of business, social, and labor force surveys. This type of research data is held within data archives or data libraries after it has been collected, so that it may be reused by future researchers. In performing their research, the detailed person-level data is aggregated into less confidential multi-dimensional tables which answer particular research questions. Portals harvest metadata (as well as publicly available data) from multiple data providers in form of RDF. To ensure high quality, the metadata must satisfy certain criteria - specified in terms of RDF constraints. After validating the metadata according to these constraints, portals offer added values to their customers, e.g., by searching over and comparing metadata of multiple providers.

By its nature, person-level data is highly confidential and access is often only permitted for qualified researchers who must apply for access. The purpose of publicly available aggregated data, on the other hand, is to get a first overview and to gain an interest in further analyses on the underlying person-level data.

Researchers typically represent their results as aggregated data in form of two-dimensional tables with only a few columns (so-called *variables* such as *sex* or *age*). The *RDF Data Cube Vocabulary (QB)*<sup>11</sup> is a W3C recommendation for representing metadata on *data cubes*, i.e. multi-dimensional aggregated data, in RDF [6]. Aggregated data is derived from person-level data by statistics on groups or aggregates such as counts, means, and frequencies. While *Disco* and *QB* provide terms for the description of data sets, both on a different level of aggregation, the *Data Catalog Vocabulary (DCAT)*<sup>12</sup> enables the representation of these data sets inside of data collections like repositories, catalogs, or archives. The relationship between data collections and their contained data sets is useful, since such collections are a typical entry point when searching for data. Aggregated data is more and more published as CSV files, allowing to perform first data calculations. In 2014, SBE and Linked Data community members developed the *Physical Data Description (PHDD)*<sup>13</sup> vocabulary to represent data in a rectangular format. The data could be either represented in records with character-separated values (CSV) or in records with fixed length.

For more detailed analyses, researchers refer to person-level data from which aggregated data is derived from, as person-level data include additional variables needed for further research. One very common example for detailed analyses on person-level data is the content-driven comparison of multiple studies. A *study* represents the process by which a data set was generated or collected. Eurostat<sup>14</sup> is the statistical office of the European Union. Its task is to provide statistics at European level that enable comparisons between countries and regions. Eurostat provides publicly available European aggregated data (downloadable as CSV files) and its metadata. This way, researchers get promising findings (in form of published tables with a few columns), e.g., about the availability of childcare services in European Union Member States by year, duration, and child age leading to subsequent research questions.

The variable *formal childcare*<sup>15</sup> (in contrast to childcare at home) captures the measured availability of childcare services in percent over the population. The present data collection refers to data on formal childcare by the variables *year*, *duration* (in hours per week), *age* of the child, and *country*. Variables are constructed out of values (of one or multiple datatypes) and/or code lists. The variable *age*, e.g., may be represented by values of the datatype *xsd:nonNegativeInteger*, or by a code list including multiple age clusters (such as '0 to 10' and '11 to 20'). To determine if variables measuring *age* - collected for different countries (*age<sub>DE</sub>*, *age<sub>UK</sub>*) - are comparable, both content-driven and technology-driven constraints are validated. Content-driven constraints ensure that the data is consistent with the intended syntax, semantics, and integrity

<sup>11</sup> <http://www.w3.org/TR/vocab-data-cube/>

<sup>12</sup> <http://www.w3.org/TR/vocab-dcat/>

<sup>13</sup> <https://github.com/linked-statistics/physical-data-description>

<sup>14</sup> <http://ec.europa.eu/eurostat>

<sup>15</sup> The data set is available at: [http://ec.europa.eu/eurostat/web/products-datasets/-/ilc\\_caindformal](http://ec.europa.eu/eurostat/web/products-datasets/-/ilc_caindformal)

of vocabularies' data models and technology-driven constraints can be generated completely automatically out of vocabularies' data models. Examples for content-driven constraints are to investigate (1) if variables are represented in a compatible way, i.e. are the variables' code lists theoretically comparable, and (2) if variables' code lists are properly structured. With technology-driven constraints, it can be validated (1) if variable definitions are available and (2) if for each code an associated category (a human-readable label) is specified.

Data providers and harvesters do not only offer metadata but also publicly available data on different level of aggregation. To ensure high data quality, they have to check provenance information and to analyze and therefore validate the data according to predefined constraints (e.g. 'are fundamental data fragments available?', and 'how does valid data look like?').

### 3 Vocabularies to Represent Metadata and Data in RDF

In this section, we describe how data in rectangular format (*PHDD*) and metadata on person-level data sets (*Disco*), aggregated data sets (*QB*), thesauri (*SKOS*), and statistical classifications (*XKOS*) are represented in RDF<sup>16</sup> and how therefore reused vocabularies are interrelated.

**Metadata on Aggregated Data.** The vocabulary *QB* represents metadata on multi-dimensional aggregate data in two files, a *qb:DataSet* and a *qb:DataStructureDefinition*. The *qb:DataStructureDefinition* contains metadata of the present data collection. Thereby, the variable *formal childcare* is modelled as *qb:measure*, since it stands for what has been measured in the data collection. The variables *year*, *duration*, *age*, and *country* are defined as *qb:dimension*. Data values, i.e., the availability of childcare services in percent over the population, are collected in a *qb:DataSet*. Each data value is represented inside a *qb:Observation* which additionally contains values for each dimension (e.g., the year in which *formal childcare* has been determined).

**Rectangular Data.** *PHDD* represents data in a rectangular format in RDF. The data could be either represented in records with character-separated values (CSV) or fixed length. Eurostat provides the two-dimensional table about *formal childcare* in form of a CSV file. The *phdd:Table* is structured by a table structure (*phdd:TableStructure*, *phdd:Delimited*). The table structure includes information about the character set (*ASCII*), the variable delimiter (*,*), the new line marker (*CRLF*), and the first line where the data starts (*2*). The table structure is related to table columns (*phdd:Column*) which are described by column descriptions (*phdd:DelimitedColumnDescription*). For the column containing the cell values in percent, the column position (*5*), the recommended data type (*xsd:nonNegativeInteger*), and the storage format (*TINYINT*) is stated. The RDFication enables further aggregations and calculations, e.g., in order to compare *formal childcare* between Northern and Southern Europe or between otherwise grouped countries.

<sup>16</sup> The complete running example in RDF is available at: <https://github.com/boschthomas/rdf-validation/tree/master/data/running-example>

**Metadata on Person-Level Data.** For a broader view of the data framework and more detailed analyses we refer to the metadata on person-level data collected for the series *EU-SILC (European Union Statistics on Income and Living Conditions)*<sup>17</sup> published by the *Microdata Information System (MISSY)*<sup>18</sup>. Where data collection is cyclic, data sets may be released as *series*, where each cycle of the data collection activity produces one or more data sets. *Missy* is an online service platform that provides systematically structured metadata for official statistics on European person-level data sets. Aggregated (qb:DataSet) and underlying person-level data sets (*disco:LogicalDataSet*) are connected by *prov:wasDerivedFrom*. The aggregated variable *formal childcare* is calculated on the basis of six person-level variables like *Education at pre-school*<sup>19</sup>. For each person-level variable detailed metadata is given (definitions, descriptions, theoretical concepts, questions variables are based on, code lists, frequencies, descriptive statistics, countries, year of data collection, and classifications) which enables researchers to replicate the results shown in the aggregated data tables from Eurostat. The vocabulary *Disco* represents metadata on person-level data in RDF. The series (*disco:StudyGroup*) *EU-SILC* contains one study (*disco:Study*) for each year (*dcterms:temporal*) of data collection. *dcterms:spatial* points to the countries for which the data has been collected. The study *EU-SILC 2011* contains eight person-level data sets (*disco:LogicalDataSet*) including person-level variables (*disco:Variable*) like the six ones needed to calculate the aggregated variable *formal childcare*.

**Organizations, Hierarchies, and Classifications.** The *Simple Knowledge Organization System (SKOS)* is reused multiple times to represent metadata on aggregated and person-level data. Variables are constructed out of values and/or (un)ordered code lists. The codes of the variable *Education at pre-school* (number of education hours per week) are modeled as *skos:Concepts* and a *skos:OrderedCollection* organizes them in a particular order within a *skos:memberList*. A variable may be associated with a theoretical concept (*skos:Concept*) and *skos:narrower* builds the hierarchy of theoretical concepts within the *skos:ConceptScheme* of a series. The variable *Education at pre-school*, e.g., is assigned to the theoretical concept *Child Care* which is the narrower concept of *Education* - one of the top concepts of the series *EU-SILC*. Controlled vocabularies (*skos:ConceptScheme*), serving as extension and reuse mechanism, organize types (*skos:Concept*) of descriptive statistics (*disco:SummaryStatistics*) like minimum, maximum, and arithmetic mean. *XKOS*<sup>20</sup> is a SKOS extension to describe formal statistical classifications like the International Standard Classification of Occupations (*ISCO*).

**Searching for (Meta)data.** *DCAT* enables to represent aggregated and person-level data sets inside of data collections like portals, repositories, catalogs, and archives which serve as typical entry points when searching for data.

<sup>17</sup> <http://www.geis.org/missy/eu/metadata/EU-SILC>

<sup>18</sup> <http://www.geis.org/missy/eu/missy-home>

<sup>19</sup> <http://www.geis.org/missy/eu/metadata/EU-SILC/2011/Cross-sectional/original#2011-Cross-sectional-RL010>

<sup>20</sup> <https://github.com/linked-statistics/xkos>

Users search for aggregated and person-level data records (*dcat:CatalogRecord*) inside data catalogs (*dcat:Catalog*). As search differs depending on the users' information need, users may only search for records' metadata (e.g., *dcterms:title*, *dcterms:description*), or may formulate more sophisticated queries on aggregated and person-level data sets (*dcat:Dataset*) or their distributions (*dcat:Distribution*) which are part of the records. Often, users search for data sets covering particular topics (*dcat:keyword*, *dcat:theme*), time periods (*dcterms:temporal*), or locations (*dcterms:spatial*), or for certain formats in which the data distribution is available (*dcterms:format*).

## 4 RDF Constraints Classification

Bosch et al. identified 74 requirements to formulate RDF constraints (e.g. *R-75*, *R-81: minimum qualified cardinality restrictions*); each of them corresponding to an RDF constraint type<sup>21</sup>[3]. We published a technical report<sup>22</sup> in which we explain each requirement (constraint type) in detail and give examples for each expressed by different constraint languages. The knowledge representation formalism *Description logics (DL)*, with its well-studied theoretical properties, provides the foundational basis for each constraint type. Therefore, this technical report contains mappings to DL to logically underpin each requirement and to determine which DL constructs are needed to express each constraint type [3].

We developed a system to classify RDF constraints and RDF constraint types. Constraints are instantiated from constraint types in order to validate both metadata and data represented by any vocabulary; thus, the proposed classification system is vocabulary-independent and therefore applicable generically. The complete set of *constraint types (CT)* encompasses two disjoint **sets of constraint types**:

1. *CT<sub>C</sub>: Content-Driven Constraint Types*
2. *CT<sub>T</sub>: Technology-Driven Constraint Types*

*CT<sub>C</sub> (Content-Driven Constraint Types)* is the set of constraints ensuring that the data is consistent with the intended syntax, semantics, and integrity of vocabularies' data models (section 4.2). Therefore, domain experts may formulate *CT<sub>C</sub>* constraints manually in prose English which is technically expressed by a constraint language afterwards. For assessing the quality of thesauri, e.g., we concentrate on the graph-based structure and apply graph- and network-analysis techniques (constraint type: *structure*). A thesaurus, e.g., should not contain many orphan concepts (concepts without any associative or hierarchical relations) lacking valuable context information for search and retrieval. *CT<sub>T</sub> (Technology-Driven Constraint Types)* is the set of constraints which can be generated completely automatically out of vocabularies' data models (section

<sup>21</sup> Constraint types and constraints are uniquely identified by alphanumeric technical identifiers like *R-71-CONDITIONAL-PROPERTIES*

<sup>22</sup> Available at: <http://arxiv.org/abs/1501.03933>

4.1). With *minimum qualified cardinality restrictions* (R-74), e.g., one can restrict that a *phdd:TableStructure* has (*phdd:column*) at least one *phdd:Column* (in DL:  $\text{TableStructure} \sqsubseteq \geq 1 \text{ column.Column}$ )<sup>23</sup>.

We determined the default **severity level** (corresponds to requirement R-158) for each constraint to indicate how serious the violation of the constraint is. We propose an extensible metric to measure the continuum of severity levels ranging from  $\mathcal{SL}_0$  to  $\mathcal{SL}_2$ . According to the constraints' default severity level the complete set of constraints ( $\mathcal{C}$ ) encompasses three disjoint **sets of constraints**:

- $\mathcal{SL}_0$ : set of constraints with severity level *informational*
- $\mathcal{SL}_1$ : set of constraints with severity level *warning*
- $\mathcal{SL}_2$ : set of constraints with severity level *error*

Violations of  $\mathcal{SL}_0$  constraints point to possible data improvements to achieve RDF representations which are ideal in terms of syntax and semantics of used vocabularies. Data not conforming to  $\mathcal{SL}_1$  and  $\mathcal{SL}_2$  constraints is syntactically or semantically not correctly represented. The difference between  $\mathcal{SL}_1$  and  $\mathcal{SL}_2$  constraints is that  $\mathcal{SL}_1$  invalid data could be whereas  $\mathcal{SL}_2$  invalid data cannot be processed further. Although, we provide default severity levels for each constraint, users should be able to specify constraints' severity levels according to their individual needs, i.e., use case specific severity levels.

We recently published a technical report<sup>24</sup> (serving as first appendix of this paper) in which we describe 212 constraints (classified as 53 constraint types) to validate rectangular data (*PHDD*) and metadata on person-level data sets (*Disco*), aggregated data sets (*QB*), thesauri (*SKOS*), and statistical classifications (*XKOS*), and therefore apply the proposed classification system to several vocabularies to represent both data and metadata [5]. In this section, we describe constraints, which are important to ensure *SBE* (meta)data quality, associate them with default severity levels, and assign them to  $\mathcal{CT}_C$  and  $\mathcal{CT}_T$  constraint types.

#### 4.1 Technology-Driven Constraint Types

$\mathcal{CT}_T$  constraints can be directly and automatically derived from explicitly stated syntax and semantics of vocabularies' data models. Thus, associated default severity levels are in most cases very strong ( $\mathcal{SL}_2$ ).

**Vocabulary.** One should not invent new or use deprecated terms (e.g. *disco:containsVariable*) of vocabularies (*vocabulary*). *Property domain* (R-25, R-26) and *range* (R-28, R-35) constraints restrict domains and ranges of properties. Only *phdd:Tables*, e.g., can have *phdd:isStructuredBy* relationships ( $\exists \text{ isStructuredBy.T} \sqsubseteq \text{Table}$ ) and *xkos:belongsTo* relationships can only point to *skos:Concepts* ( $\top \sqsubseteq \forall \text{ belongsTo.Concept}$ ). A *universal quantification* (R-91) contains all those individuals that are connected by a property only to individuals/literals of particular classes or data ranges. Only *dcat:Catalogs*, e.g., can have *dcat:dataset* relationships to *dcat:Datasets*

<sup>23</sup> For simplicity reasons, we do not use namespace prefixes in DL statements.

<sup>24</sup> Available at: <http://arxiv.org/abs/1504.04479>



(`Catalog`  $\sqsubseteq \forall$  `dataset.Dataset`). Out-dated classes and properties of previous vocabulary versions can be marked as deprecated. The constraint types *context-specific valid classes and properties* (*R-209*; *R-210*) can be used to specify which classes and properties are valid in which context - here a given vocabulary version. Many properties are not necessarily required but *recommended* within a particular context (*R-72*). The property *skos:notation*, e.g., is not mandatory for *disco:Variables*, but recommended to represent variable names. *R-223* serves to make sure that all literal values are valid with regard to their datatypes. Thus, all date values (e.g. *disco:startDate*, *disco:endDate*, *dcterms:date*) must be of the datatype *xsd:date* and *xsd:nonNegativeInteger* values (e.g. *disco:frequency*) do not have to be negative.

**Existential Quantifications.** *Existential quantifications* (*R-86*) enforce that instances of given classes must have some property relation to individuals/literals of certain types. If a study, e.g., does not contain any data set ( $\mathcal{SL}_2$ ), the actual description of the data is missing which may indicate that it is very hard or unlikely to get access to the data. If metadata on data files, including the actual data, is missing (especially case and variable quantities;  $\mathcal{SL}_1$ ), the description of the data sets and the study is not sufficient. Case quantity measures how many cases are collected for a study. High case and variable quantities are indicators for high statistical quality and comprehensiveness of the underlying study ( $\mathcal{SL}_1$ ). Variables should have a relation to a theoretical concept ( $\mathcal{SL}_0$ ). The variable *Education at pre-school*, e.g., is associated with the theoretical concept *Child Care*. The constraint's severity level is weak, as in most cases research can be continued without having information about theoretical concepts.

**Cardinality Restrictions on Properties and Language Tags.** *Minimum/maximum/exact qualified cardinality restrictions* (*R-74*, *R-75*, *R-76*) contain all those individuals that are connected by a property to at least/at most/exactly  $n$  different individuals/literals of particular classes or data ranges. A *phdd:TableStructure*, e.g., has (*phdd:column*) at least one *phdd:Column* (`TableStructure`  $\sqsubseteq \geq 1$  `column.Column`), a *disco:Variable* has at most one *disco:concept* relationship to a theoretical concept (*skos:Concept*) (`Variable`  $\sqsubseteq \leq 1$  `concept.Concept`), and a *qb:DataSet* is structured by (*qb:structure*) exactly one *qb:DataStructureDefinition* (`DataSet`  $\sqsubseteq \geq 1$  `structure.DataStructureDefinition`  $\sqcap \leq 1$  `structure.DataStructureDefinition`). For data properties, it may be desirable to restrict that values of predefined languages must be present for determined number of times (*R-48*, *R-49*): (1) Some controlled vocabularies contain literals in natural language, but without information what language has actually been used. (2) Language tags must conform to language standards. (3) Some thesaurus concepts are labeled in only one, others in multiple languages. It may be desirable to have each concept labeled in each of the languages that are also used on the other concepts. Although not always possible, incompleteness of language coverage for some concepts may indicate shortcomings of thesauri [9].

**Disjointness and Allowed Values.** All properties, not having the same domain and range types, are defined to be pairwise disjoint (*R-9: disjoint properties*), i.e., no individual  $x$  can be connected to an individual/literal  $y$  by dis-

joint properties like *phdd:isStructuredBy* and *phdd:column* (*isStructuredBy*  $\sqsubseteq \neg \text{column}$ ). All *XKOS* classes are pairwise disjoint (*R-7: disjoint classes*; e.g. *ClassificationLevel*  $\sqcap$  *ConceptAssociation*  $\sqsubseteq \perp$ ), i.e., individuals cannot be instances of multiple disjoint classes. It is a common requirement to narrow down the value space of properties by an exhaustive enumeration of valid values. *Allowed values* (*R-30, R-37*) for properties can be IRIs (matching one or multiple patterns), any literals, allowed literals (e.g. 'red' 'blue' 'green'), and typed literals of one or multiple type(s) (e.g. *xsd:string*). *disco:CategoryStatistics*, e.g., can only have *disco:computationBase* relationships to the values *valid* and *invalid* of the datatype *rdf:langString* (*CategoryStatistics*  $\equiv \forall \text{computationBase}.\{\text{valid}, \text{invalid}\} \sqcap \text{langString}$ ).

**Validation and Reasoning.** Some constraint types enable performing reasoning prior to validation which may resolve or cause constraint violations. With *subsumption* (*R-100*), one can state that *xkos:ClassificationLevel* is a sub-class of *skos:Collection*, i.e., each *xkos:ClassificationLevel* must also be part of the *skos:Collection* class extension (*ClassificationLevel*  $\sqsubseteq$  *Collection*). With *sub properties* (*R-54, R-64*), one can state that *disco:fundedBy* is a sub-property of *dcterms:contributor* - i.e., if a study is funded by an organization, then this organization contributed to this study (*fundedBy*  $\sqsubseteq$  *contributor*). *Default values* (*R-31, R-38*) for objects/literals of given properties are inferred automatically when properties are not present in the data. The value *true* for the property *disco:isPublic* indicates that a *disco:LogicalDataSet* can be accessed by anyone. Per default, however, access to data sets should be restricted (*false*). Validation should *exploit sub-super relations* in vocabularies (*R-224*). If *dcterms:coverage* and one of its sub-properties (*dcterms:spatial*, *dcterms:temporal*) are present, it is checked that *dcterms:coverage* is not redundant with its sub-properties. This validation can indicate when the data is verbose (redundant) or expressed at a too general level and could thus be improved.

## 4.2 Content-Driven Constraint Types

In this sub-chapter, we assign default severity levels to and describe constraints of diverse  $\mathcal{CT}_C$  constraint types to ensure that the data is consistent with the intended syntax, semantics, and integrity of vocabularies' data models.

**Data Model Consistency.** The purpose of some constraints is to ensure the integrity of the data according to intended data model semantics (*data model consistency*). Every *qb:Observation*, e.g., must have a value for each dimension declared in its *qb:DataStructureDefinition* ( $\mathcal{SL}_2$ ) and no two *qb:Observations* in the same *qb:DataSet* can have the same value for all dimensions ( $\mathcal{SL}_1$ ). If a *qb:DataSet* *D* has a *qb:Slice* *S*, and *S* has an *qb:Observation* *O*, then the *qb:DataSet* corresponding to *O* must be *D* ( $\mathcal{SL}_1$ ). Relative frequencies of variable codes are calculated correctly, if the cumulative percentage (*disco:cumulativePercentage*) of a given code exactly matches the cumulative percentage of the previous code plus the percentage value (*disco:percentage*) of the current code ( $\mathcal{SL}_2$ ).

**Structure and Ordering.** For assessing the quality of *SKOS* vocabularies, we concentrate on graph-based structures and apply graph- and network-analysis

techniques (*structure*) like (1) a thesaurus should provide entry points (top concepts) to the data to provide efficient access and guidance for human users, (2) concepts, internal to the tree, should not be indicated as top concepts, and (3) a thesaurus should not contain many orphan concepts (concepts without any associative or hierarchical relations) lacking valuable context information for retrieval. Objects/literals can be declared to be *ordered* (R-121, R-217) for given properties. Variables, questions, and codes, e.g., are typically organized in a particular order. If codes (*skos:Concept*) should be ordered, they must be members (*skos:memberList*) in an ordered collection (*skos:OrderedCollection*), the variable's code list.

**Comparison.** A very common research question is to compare variables of multiple studies or countries (*comparison*). To compare variables, (1) variables and (2) variable definitions must be present, (3) code lists must be structured properly, (4) for each code an associated category must be specified, and (5) code lists must either be identical or at least similar. If a researcher wants to get a first overview over comparable variables (use case 1), covering the first three constraints may be sufficient. Thus, the severity level of the first three constraints is stronger ( $\mathcal{SL}_2$ ) than for the last two constraints ( $\mathcal{SL}_1$  and  $\mathcal{SL}_0$ ). If the intention of the researcher is to perform more detailed comparisons (use case 2), however, the violation of the remaining two constraints is getting more serious.

**Unique Identification.** It is often useful to declare a given (data) property as the *primary key* (R-226) of a class, so that a system can enforce uniqueness and also automatically build URIs from user inputs and imported data. In *Disco*, resources are uniquely identified by the property *adms:identifier*, which is therefore inverse-functional (**func** **identifier**<sup>-</sup>), i.e. for each *rdfs:Resource* *x*, there can be at most one distinct *rdfs:Resource* *y* such that *y* is connected by *adms:identifier*<sup>-</sup> to *x* ( $\mathcal{SL}_2$ ). Keys, however, are even more general than *inverse-functional properties* (R-58), as a key can be a data, an object property, or a chain of properties [10]. For this generalization purposes, as there are different sorts of key, and as keys can lead to undecidability, DL is extended with *key boxes* and the construct *keyfor* (**identifier** **keyfor** **Resource**) [8]. OWL 2 *HasKey* implements *keyfor* ( $\mathcal{SL}_2$ ) and thus can be used to identify resources uniquely, to merge resources with identical key property values, and to recognize constraint violations.

**Membership in Controlled Vocabularies.** In many cases, resources must be *members of controlled vocabularies* (R-32). If a dimension property, e.g., has a *qb:codeList*, then the value of the dimension property on every *qb:Observation* must be in the code list ( $\mathcal{SL}_2$ ). Summary statistics types like minimum, maximum, and arithmetic mean are maintained within a controlled vocabulary. Thus, summary statistics can only have *disco:summaryStatisticType* relationships to *skos:Concepts* which must be members of the controlled vocabulary *ddicv:SummaryStatisticType*, a *skos:ConceptScheme* ( $\mathcal{SL}_2$ ).

**Constraints on Properties.** It is useful to declare properties to be *conditional* (R-71), i.e., if particular properties exist (or do not exist), then other prop-

erties must also be present (or absent). To get an overview over a series/study either an abstract, a title, an alternative title, or links to external descriptions should be provided. If an abstract and an external description are absent, however, a title or an alternative title should be given ( $\mathcal{SL}_1$ ). In case a variable is represented in form of a code list, codes may be associated with categories, i.e., human-readable labels ( $\mathcal{SL}_0$ ). The variable *Education at pre-school*, e.g., is represented as ordered code list without any categories. If a *skos:Concept* represents a code (having *skos:notation* and *skos:prefLabel* properties), then the property *disco:is Valid* has to be stated indicating if the code stands for valid (*true*) or missing (*false*) cases ( $\mathcal{SL}_2$ ). Constraints of type *context-specific exclusive or of property groups* (*R-11*) restrict individuals of given classes to have properties defined within exactly one of multiple property groups. *skos:Concepts* can have either *skos:definition* (when interpreted as theoretical concepts) or *skos:notation* and *skos:prefLabel* properties (when interpreted as codes/categories), but not both ( $\mathcal{SL}_2$ ). For datatype properties, it should be possible to declare frequently needed *facets* (*R-46*) to validate input against simple conditions including min/max values, regular expressions, and string length. The abstract of series/studies, e.g., should have a minimum length ( $\mathcal{SL}_1$ ).

**Constraints on Literals.** *disco:percentage* stands for the number of cases of a given code in relation to the total number of cases for a particular variable. Percentage values are only valid when they are within the *literal range* of 0 and 100 (*R-45*;  $\mathcal{SL}_2$ ). *Mathematical Operations* (*R-41*, *R-42*; e.g. date calculations and statistical computations like average, mean, and sum) are performed to ensure the integrity of data models. The sum of percentage values of all variable codes, e.g., must exactly be 100 ( $\mathcal{SL}_2$ ) and the minimum absolute frequency of all variable codes do not have to be greater than the maximum ( $\mathcal{SL}_2$ ). Depending on property datatypes, two different literal values have a specific ordering with respect to an operator like *<* (*R-43: literal value comparison*). Start dates (*disco:startDate*), e.g., must be before (*<*) end dates (*disco:endDate*) ( $\mathcal{SL}_2$ ).

## 5 Implementation

SPARQL is generally seen as the method of choice to validate RDF data according to certain constraints. We use *SPIN*, a SPARQL-based way to formulate and check constraints, as basis to develop a validation environment (available at <http://purl.org/net/rdfval-demo>)<sup>25</sup> to validate RDF data according to constraints expressed by arbitrary constraint languages like Shape Expressions, Resource Shapes, and the Web Ontology Language<sup>26</sup> [2]. The *RDF Validator* also validates RDF data to ensure correct syntax, semantics, and integrity of diverse vocabularies such as *Disco*, *QB*, *PHDD*, *SKOS*, and *XKOS*. Although accessible within our validation tool, we provide all implemented constraints<sup>27</sup> in

<sup>25</sup> Source code downloadable at: <https://github.com/boschthomas/rdf-validator>

<sup>26</sup> SPIN mappings available at: <https://github.com/boschthomas/rdf-validation/tree/master/SPIN>

<sup>27</sup> <https://github.com/boschthomas/rdf-validation/tree/master/constraints>

form of SPARQL CONSTRUCT queries. For the subsequent evaluation, we implemented 212 constraints on *Disco*, *QB*, *SKOS*, *XKOS*, and *PHDD* data sets. The SPIN engine checks for each resource if it satisfies all constraints, which are associated with its assigned classes, and generates a result RDF graph containing information about all constraint violations. There is one SPIN construct template for each constraint type and vocabulary-specific constraint<sup>28</sup>. A SPIN construct template contains a SPARQL CONSTRUCT query which generates constraint violation triples indicating the subject and the properties causing constraint violations, and the reason why constraint violations have been raised. A SPIN construct template creates constraint violation triples if all triple patterns within the SPARQL WHERE clause match. *Missy*<sup>29</sup> provides comprehensive Linked Data services like diverse RDF exports of person-level metadata conforming to the *Disco* vocabulary in form of multiple concrete syntaxes.

## 6 Evaluation

### 6.1 Evaluation Setup

First, several SBE domain experts of the vocabularies *Disco*, *QB*, *SKOS*, *XKOS*, and *PHDD* evaluated the correctness (i.e., the gold standard) of all  $\mathcal{CT}_C$  and  $\mathcal{CT}_T$  constraints and therefore the generic applicability of the developed classification system of constraint types and constraints. Second, we exhaustively evaluated the metadata quality of large real world aggregated (*QB*), person-level (*Disco*), and thesauri (*SKOS*) data sets by means of both  $\mathcal{C}_C$  and  $\mathcal{C}_T$  constraints of the majority of the constraint types. We validated 9,990 / 3,775,983,610 (*QB*), 4,178 / 477,737,281 (*SKOS*), and 1,526 / 9,673,055 (*Disco*) data sets / triples using the *RDF Validator* in batch mode. That are more than 4.2 billion triples and 15 thousand data sets. We validated, i.a., (1) *QB* data sets published by the *Australian Bureau of Statistics (ABS)*, the *European Central Bank (ECB)*, and the *Organisation for Economic Co-operation and Development (OECD)*, (2) *SKOS* thesauri like the *AGROVOC Multilingual agricultural thesaurus*, the *STW Thesaurus for Economics*, and the *Thesaurus for the Social Sciences (TheSoz)*, and (3) *Disco* data sets provided by the *Microdata Information System (Missy)*, the *DwB Discovery Portal*, the *Danish Data Archive (DDA)*, and the *Swedish National Data Service (SND)*.

We recently published a technical report<sup>30</sup> (serving as second appendix of this paper) in which we describe the comprehensive evaluation in detail [4]. As we evaluated nearly 10 thousand *QB* data sets, we published the evaluation results for each data set in form of one document per SPARQL endpoint<sup>31</sup>. Table 1 shows the evaluation results.

<sup>28</sup> For a comprehensive description of the *RDF Validator*, we refer to [2]

<sup>29</sup> <http://www.gesis.org/missy/eu/missy-home>

<sup>30</sup> Available at: <http://arxiv.org/abs/1504.04478>

<sup>31</sup> Available at: <https://github.com/boschthomas/rdf-validation/tree/master/evaluation/data-sets/data-cube>

## 6.2 Evaluation Results and Discussion

Criteria	<i>Disco</i>	<i>QB</i>	<i>SKOS</i>	Total
<i>Triples</i>	9,673,055	3,775,983,610	477,737,281	4,263,393,946
<i>Data Sets</i>	1,526	9,990	4,178	15,694
<i>CV</i>	3,545,703	45,635,846	5,540,988	54,722,537
<i>CV</i> ( $\mathcal{SL}_0$ )	2,437,922 ( <b>68.8%</b> )	0 (0%)	2,281,740 (41.2%)	4,719,662 (8.6%)
<i>CV</i> ( $\mathcal{SL}_1$ )	473,574 (13.4%)	45,520,613 ( <b>99.75%</b> )	3,259,248 ( <b>58.8%</b> )	49,253,435 ( <b>90%</b> )
<i>CV</i> ( $\mathcal{SL}_2$ )	634,207 (17.9%)	115,233 (0.25%)	0 (0%)	749,440 (1.4%)
<i>CT</i>	52	20	14	53
<i>CT</i> ( $\mathcal{C}_C$ )	30 ( <b>57.7%</b> )	5 (25%)	5 (35.7%)	30 ( <b>56.6%</b> )
<i>CT</i> ( $\mathcal{C}_T$ )	22 (42.3%)	15 ( <b>75%</b> )	9 ( <b>64.3%</b> )	23 (43.4%)
<i>C</i>	142	35	35	212
<i>C</i> ( $\mathcal{C}_C$ )	72 ( <b>50.7%</b> )	16 (45.7%)	21 ( <b>60%</b> )	109 ( <b>51.4%</b> )
<i>C</i> ( $\mathcal{C}_T$ )	70 ( <b>49.3%</b> )	19 ( <b>54.3%</b> )	14 (40%)	103 ( <b>48.6%</b> )
<i>C</i> ( $\mathcal{SL}_0$ )	75 ( <b>52.8%</b> )	4 (11.4%)	21 ( <b>60%</b> )	100 ( <b>47.2%</b> )
<i>C</i> ( $\mathcal{SL}_1$ )	9 (6.3%)	3 (8.6%)	5 (14.3%)	17 (8%)
<i>C</i> ( $\mathcal{SL}_2$ )	58 (40.8%)	28 ( <b>80%</b> )	9 (25.7%)	95 ( <b>44.8%</b> )

*C* (constraints), *CT* (constraint types), *CV* (constraint violations)

Table 1: Evaluation

We identified 142 *Disco* constraints ( $\mathcal{C}_C$  and  $\mathcal{C}_T$  constraints to the same extend) assigned to 52 distinct constraint types and implemented 77 of them to actually validate person-level data sets. For *QB*, we specified more  $\mathcal{C}_T$  (54%) than  $\mathcal{C}_C$  constraints; for *SKOS*, however, more  $\mathcal{C}_C$  constraints (60%). We instantiated more  $\mathcal{C}_C$  (58%) than  $\mathcal{C}_T$  constraint types to define *Disco* constraints; for *QB* (75%) and *SKOS* (64%), on the other side, more  $\mathcal{C}_T$  constraint types. In total, we used 53 of overall 82 distinct constraint types (57% of them are  $\mathcal{C}_C$  constraint types) to define 212 constraints (equally  $\mathcal{C}_C$  and  $\mathcal{C}_T$  constraints).

For *Disco* and *SKOS*, more than the half of the constraints are associated with the weakest severity level  $\mathcal{SL}_0$ . Within the context of *QB*, 80% of the constraints are classified as the most serious ones ( $\mathcal{SL}_2$ ). All in all, there are a little bit more  $\mathcal{SL}_0$  then  $\mathcal{SL}_2$  constraints, whereas  $\mathcal{SL}_1$  constraints are negligible. *Existential quantifications* (32.4%, *Disco*), *data model consistency* (31.4%, *QB*), and *structure* (28.6%, *SKOS*) are the constraint types the most constraints are instantiated from. By validating *QB* data sets, we got the most constraint violations (more than 45 millions), followed by *SKOS* and *Disco* (with more than 5.5 and 3.5 millions) - consequently, almost 55 million constraint violations were raised during the evaluation which could be used to enhance the metadata quality of these data sets. Close to 70% of all *Disco* constraint violations are caused by violating  $\mathcal{SL}_0$  constraints. For *QB* (nearly 100%) and *SKOS* (almost 60%), the majority of the raised constraint violations are classified to be more serious ( $\mathcal{SL}_1$ ). 80% of all *QB* constraints are  $\mathcal{SL}_2$  constraints leading to less than 1% of all *QB* constraint violations. Al-

together, exactly 90% of the constraint violations are assigned to the severity level  $\mathcal{SL}_1$ . These findings are surprising as only 8% of all defined constraints are  $\mathcal{SL}_1$  constraints. The constraints responsible for the largest numbers of constraint violations are *DISCO-C-LABELING-AND-DOCUMENTATION-06* and *DISCO-C-COMPARISON-VARIABLES-02* (both 547,916) (*Disco*), *DATA-CUBE-C-DATA-MODEL-CONSISTENCY-05* (45,514,102) (*QB*), and *SKOS-C-LANGUAGE-TAG-CARDINALITY-01* (2,508,903) (*SKOS*). We refer to the technical reports<sup>24 30</sup> to get details about constraints on and the evaluation of *XKOS* and *PHDD* data sets.

## 7 Related Work

For data archives, research institutes, and data libraries, RDF validation according to predefined constraints is a much sought-after feature, particularly as this is taken for granted in the XML world. DDI-XML documents, e.g., are validated against diverse XSDs<sup>4</sup>. As certain constraints cannot be formulated and validated by XSDs, so-called secondary-level validation tools like *Schematron*<sup>32</sup> have been introduced to overcome the limitations of XML validation. *Schematron* generates validation rules and validates XML documents according to them. With RDF validation, one can overcome drawbacks when validating XML documents<sup>33</sup>. It cannot be validated, e.g., if each code of a variable's code list is associated with a category (*R-86*). Additionally, it cannot be validated that if an element has a specific value, then certain child elements must be present (*R-71*). A comprehensive comparison of XML and RDF validation, however, is not within the scope of this paper.

A well-formed *RDF Data Cube* is an a RDF graph describing one or more instances of *qb:DataSet* for which each of the 22 integrity constraints<sup>34</sup>, defined within the *QB* specification, passes. Each integrity constraint is expressed as narrative prose and, where possible, a SPARQL ASK query or query template. If the ASK query is applied to an RDF graph then it will return true if that graph contains one or more *QB* instances which violate the corresponding constraint [7]. Mader, Haslhofer, and Isaac investigated how to support taxonomists in improving SKOS vocabularies by pointing out quality issues that go beyond the integrity constraints defined in the SKOS specification [9].

## 8 Conclusion and Future Work

In this paper, we showed in form of a complete real world running example how to represent metadata on person-level data (*Disco*), metadata on aggregated data (*QB*), and data on both aggregation levels in a rectangular format (*PHDD*)

<sup>32</sup> <https://msdn.microsoft.com/en-us/library/aa468554.aspx>

<sup>33</sup> [http://www.xmlmind.com/xmleditor/\\_distrib/doc/xmltool/xsd\\_structure\\_limitations.html](http://www.xmlmind.com/xmleditor/_distrib/doc/xmltool/xsd_structure_limitations.html)

<sup>34</sup> <http://www.w3.org/TR/vocab-data-cube/#wf>

**Kai:** @KAI: Dieses Kapitel muss ich noch an die Änderungen im Paper anpassen!

in RDF and how therefore used vocabularies are interrelated (**contribution 1**, section 3). We explained why RDF validation is important in this context and how metadata on person-level data, aggregated data, thesauri, and statistical classifications as well as data on both aggregation levels is validated against constraints to ensure high (meta)data quality<sup>35</sup> (**contribution 2**, section ??). We distinguish two validation types: (1) *Content-Driven Validation*  $\mathcal{C}_C$  contains the set of constraints ensuring that the data is consistent with the intended syntax, semantics, and integrity of data models (section 4.2). (2) *Technology-Driven Validation*  $\mathcal{C}_T$  includes the set of constraints which can be generated automatically out of data models, such as cardinality restrictions, universal and existential quantifications, domains, and ranges (section 4.1). We determined the default *severity level* for each constraint to indicate how serious the violation of the constraint is and propose an extensible metric to measure the continuum of severity levels.

We implemented a validation environment (available at <http://purl.org/net/rdfval-demo>) to validate RDF data according to constraints expressed my arbitrary constraint languages and to ensure correct syntax, semantics, and integrity of diverse vocabularies such as *Disco*, *QB*, *PHDD*, *SKOS*, and *XKOS* (section 5). We exhaustively evaluated the metadata quality of large real world aggregated (*QB*), person-level (*Disco*), and thesauri (*SKOS*) data sets by means of 212  $\mathcal{C}_C$  and  $\mathcal{C}_T$  constraints of the majority of the constraint types. We validated more than 4.2 billion triples and 15 thousand data sets<sup>36</sup> (section 6).

## 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.
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. Thomas Bosch, Benjamin Zapolko, Joachim Wackerow, and Kai Eckert. An evaluation of metadata and data quality on person-level, aggregated, thesauri, statistical classifications, and rectangular data sets. 2015.
5. Thomas Bosch, Benjamin Zapolko, Joachim Wackerow, and Kai Eckert. Rdf constraints to validate rectangular data and metadata on person-level data, aggregated data, thesauri, and statistical classifications. 2015.
6. Richard Cyganiak, Simon Field, Arofan Gregory, Wolfgang Halb, and Jeni Tennison. Semantic statistics: Bringing together sdmx and scovo. In Christian Bizer, Tom Heath, Tim Berners-Lee, and Michael Hausenblas, editors, *Proceedings of the*

<sup>35</sup> The first appendix of this paper describing each constraint in detail is available at: <http://arxiv.org/abs/1504.04479> [5]

<sup>36</sup> The second appendix of this paper describing the evaluation in detail is available at: <http://arxiv.org/abs/1504.04478> [4].



*WWW 2010 Workshop on Linked Data on the Web*, volume 628 of *CEUR Workshop Proceedings*, 2010.

7. Richard Cyganiak and Dave Reynolds. The rdf data cube vocabulary. W3C recommendation, W3C, January 2014.
8. 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.
9. Christian Mader, Bernhard Haslhofer, and Antoine Isaac. Finding quality issues in skos vocabularies. In *Proceedings of the Second International Conference on Theory and Practice of Digital Libraries*, TPD L’12, pages 222–233, Berlin, Heidelberg, 2012. Springer-Verlag.
10. Michael Schneider. OWL 2 Web Ontology Language RDF-Based Semantics. W3C recommendation, W3C, October 2009.
11. Mary Vardigan, Pascal Heus, and Wendy Thomas. Data documentation initiative: Towards a standard for the social sciences. *International Journal of Digital Curation*, 3(1):107–113, 2008.