

# RDF Constraints to Validate

## Rectangular Data and Metadata on Person-Level Data, Aggregated Data, Thesauri, and Statistical Classifications

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.** To ensure high quality of and trust in both metadata and data, their representation in RDF must satisfy certain criteria - specified in terms of RDF constraints. From 2012 to 2015 together with other Linked Data community members and experts from the social, behavioural, and economic sciences (*SBE*), we developed diverse vocabularies to represent *SBE* metadata and rectangular data in RDF. The *DDI-RDF Discovery Vocabulary (Disco)* is designed to support the dissemination, management, and reuse of person-level data, i.e., data about individuals, households, and businesses, collected in form of responses to studies and archived for research purposes. The *RDF Data Cube Vocabulary (Data Cube)* is a W3C recommendation for expressing *data cubes*, i.e. multi-dimensional aggregate data. *Physical Data Description (PHDD)* is a vocabulary to model data in rectangular format. The data could either be represented in records with character-separated values (*CSV*) or fixed length. The *Simple Knowledge Organization System (SKOS)* is a vocabulary to build knowledge organization systems such as thesauri, classification schemes, and taxonomies. *XKOS* is a SKOS extension to describe formal statistical classifications.

In this paper, we describe RDF constraints to validate metadata on person-level data (*Disco*), aggregated data (*Data Cube*), thesauri (*SKOS*), and statistical classifications (*XKOS*) and to validate rectangular data (*PHDD*). We assign RDF constraints to RDF constraint types either corresponding to RDF validation requirements or to data model specific constraint types. This technical report is updated continuously as modifying, adding, and deleting constraints remains ongoing work.

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

## 1 Introduction

Bosch and Eckert initiated a comprehensive database<sup>3</sup> on RDF validation requirements to collect case studies, use cases, and requirements [1]. It is contin-

---

<sup>3</sup> Publicly available at <http://purl.org/net/rdf-validation>.

uously updated and used to evaluate and to compare various existing solutions for RDF constraint formulation and validation. Requirements are classified to provide a high-level view on different solutions and to facilitate a better understanding of the problem domain. Bosch et al. identified in total 76 requirements to formulate RDF constraints; each of them corresponding to a constraint type. They recently published a technical report<sup>4</sup> in which they explain each requirement (constraint type) in detail and give examples for each (represented by listed constraint languages) [2]. Where appropriate constraint types are related to complementary requirements of the RDF validation requirements database.

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 given data models. (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. 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  $\mathcal{SL}_0$  (informational) via  $\mathcal{SL}_1$  (warning) to  $\mathcal{SL}_2$  (error). Although we provide default severity levels for each constraint, users should be able to specify severity levels of constraints they need to validate for their individual use cases, i.e., users should be able to define use case specific severity levels for constraints.

The following two sections encompass RDF constraints classified as constraint types either corresponding to RDF validation requirements (section 2) or data model specific constraint types (section 3). Constraint types marked with an asterisk (\*) can be used as OWL 2 axioms. Thus, reasoners may be used prior to the validation to infer implicit triples which may resolve or cause constraint violations.  $\mathcal{C}_T$  constraint types are marked with the superscript<sup>B</sup>. As  $\mathcal{C}_T$  and  $\mathcal{C}_C$  are disjoint sets,  $\mathcal{C}_C$  constraint types are not explicitly marked.

## 2 RDF Constraint Types Corresponding to RDF Validation Requirements

### 2.1 Subsumption\*<sup>B</sup>

A *subclass axiom*<sup>5</sup> (*concept inclusion* in DL) states that the class  $C1$  is a subclass of the class  $C2$  -  $C1$  is more specific than  $C2$ , i.e. each resource of the class  $C1$  must also be part of the class extension of  $C2$ .

- **DISCO-C-SUBSUMPTION-01:** All *disco:Universes* must also be *skos:Concepts* (Universe  $\sqsubseteq$  Concept).
  - severity level: ERROR

<sup>4</sup> Available at: <http://arxiv.org/abs/1501.03933>

<sup>5</sup> *R-100-SUBSUMPTION*

## 2.2 Class Equivalence\*B

*Class Equivalence*<sup>6</sup> asserts that two concepts have the same instances. While synonyms are an obvious example of equivalent concepts, in practice one more often uses concept equivalence to give a name to complex expressions [5]. Concept equivalence is indeed subsumption from left and right ( $A \sqsubseteq B$  and  $B \sqsubseteq A$  implies  $A \equiv B$ ).

- **DISCO-C-CLASS-EQUIVALENCE-01:** All *sio:SIO\_000367* resources must also be *disco:Variables* (**Variable**  $\equiv$  **SIO\_000367**). The SemanticScience Integrated Ontology (SIO)<sup>7</sup> provides a simple, integrated ontology of types and relations for rich description of objects, processes and their attributes. *sio:SIO\_000367* is a variable defined as a value that may change within the scope of a given problem or set of operations. Thus, *sio:SIO\_000367* is equivalent to *disco:Variable*.
  - severity level: INFO

## 2.3 Sub Properties\*B

*Sub Properties*<sup>8</sup> state that the property *P1* is a sub property of the property *P2* - that is, if an individual *x* is connected by *P1* to an individual or a literal *y*, then *x* is also connected by *P2* to *y*.

- **DISCO-C-SUB-PROPERTIES-01:** If an individual *x* is connected by *disco:fundedBy* to an individual *y*, then *x* is also connected by *dterms:contributor* to *y* (**fundedBy**  $\sqsubseteq$  **contributor**).
  - severity level: ERROR

## 2.4 Property Domains\*B

*Property Domains*<sup>9</sup> (*domain restrictions on roles* in DL) restrict the domain of object and data properties. The purpose is to declare that a given property is associated with a class. In OO terms this is the declaration of a member, field, attribute or association.  $\exists R.\top \sqsubseteq C$  is the object property restriction where *R* is the object property (role) whose domain is restricted to concept *C*.

- **DISCO-C-PROPERTY-DOMAIN-01:** *Property Domain* constraints are defined for each *Disco* object and data property. Only *disco:Questions*, e.g., can have *disco:responseDomain* relationships ( $\exists$  **responseDomain**. $\top \sqsubseteq$  **Question**).
  - Severity level: ERROR

<sup>6</sup> *R-3-EQUIVALENT-CLASSES*

<sup>7</sup> <https://code.google.com/p/semanticscience/wiki/SIO>

<sup>8</sup> *R-54-SUB-OBJECT-PROPERTIES*, *R-64-SUB-DATA-PROPERTIES*

<sup>9</sup> *R-25-OBJECT-PROPERTY-DOMAIN*, *R-26-DATA-PROPERTY-DOMAIN*

- **DATA-CUBE-C-PROPERTY-DOMAIN-01:** *Property Domain* constraints are defined for each *Data Cube* object and data property. Only *qb:Observations*, e.g., can have *qb:dataset* relationships ( $\exists \text{ dataset.}\top \sqsubseteq \text{Observation}$ ).
  - Severity level: ERROR
- **DCAT-C-PROPERTY-DOMAIN-01:** *Property Domain* constraints are defined for each *DCAT* object and data property. Only *dcat:Catalogs*, e.g., can have *dcat:dataset* relationships ( $\exists \text{ dataset.}\top \sqsubseteq \text{Catalog}$ ).
  - Severity level: ERROR
- **PHDD-C-PROPERTY-DOMAIN-01:** *Property Domain* constraints are defined for each *PHDD* object and data property. Only *phdd:Tables*, e.g., can have *phdd:isStructuredBy* relationships ( $\exists \text{ isStructuredBy.}\top \sqsubseteq \text{Table}$ ).
  - Severity level: ERROR
- **SKOS-C-PROPERTY-DOMAIN-01:** *Property Domain* constraints are defined for each *SKOS* object and data property. Only *skos:ConceptSchemes*, e.g., can have *skos:hasTopConcept* relationships ( $\exists \text{ hasTopConcept.}\top \sqsubseteq \text{ConceptScheme}$ ).
  - Severity level: ERROR
- **XKOS-C-PROPERTY-DOMAIN-01:** *Property Domain* constraints are defined for each *XKOS* object and data property.
  - Severity level: ERROR

## 2.5 Property Ranges\*<sup>B</sup>

*Property Ranges*<sup>10</sup> (*range restrictions on roles* in DL) restrict the range of object and data properties.  $\top \sqsubseteq \forall R.C$  is the range restriction to the object property *R* (restricted by the concept *C*).

- **DISCO-C-PROPERTY-RANGES-01:** *Property Range* constraints are defined for each *Disco* object and data property. *disco:caseQuantity* relationships, e.g., can only point to literals of the datatype *xsd:nonNegativeInteger* ( $\top \sqsubseteq \forall \text{ caseQuantity.nonNegativeInteger}$ ).
  - Severity level: ERROR
- **DATA-CUBE-C-PROPERTY-RANGES-01:** *Property Range* constraints are defined for each *Data Cube* object and data property. *qb:order* relationships, e.g., can only point to literals of the datatype *xsd:string* ( $\top \sqsubseteq \forall \text{ order.string}$ ).
  - Severity level: ERROR
- **DCAT-C-PROPERTY-RANGES-01:** *Property Range* constraints are defined for each *DCAT* object and data property. *dcat:bytes* relationships, e.g., can only point to literals of the datatype *xsd:integer* ( $\top \sqsubseteq \forall \text{ bytes.integer}$ ).

<sup>10</sup> R-28-OBJECT-PROPERTY-RANGE, R-35-DATA-PROPERTY-RANGE

- Severity level: ERROR
- **PHDD-C-PROPERTY-RANGES-01:** *Property Range* constraints are defined for each *PHDD* object and data property. *phdd:caseQuantity* relationships, e.g., can only point to literals of the datatype *xsd:nonNegativeInteger* ( $\top \sqsubseteq \forall \text{ caseQuantity.nonNegativeInteger}$ ).
  - Severity level: ERROR
- **SKOS-C-PROPERTY-RANGES-01:** *Property Range* constraints are defined for each *SKOS* object and data property.
  - Severity level: ERROR
- **XKOS-C-PROPERTY-RANGES-01:** *Property Range* constraints are defined for each *XKOS* object and data property. *xkos:belongsTo* relationships, e.g., can only point to instances of the class *skos:Concept* ( $\top \sqsubseteq \forall \text{ belongsTo.Concept}$ ).
  - Severity level: ERROR

## 2.6 Inverse Object Properties\*<sup>B</sup>

In many cases, properties are used bi-directionally and then accessed in the inverse direction, e.g.  $\text{parent} \equiv \text{child}^-$ . There should be a way to declare value type, cardinality etc of those inverse relations without having to declare a new property URI. The object property *OP1* is an inverse<sup>11</sup> of the object property *OP2*. Thus, if an individual *x* is connected by *OP1* to an individual *y*, then *y* is also connected by *OP2* to *x*, and vice versa.

- **DISCO-C-INVERSE-OBJECT-PROPERTIES-01:** *disco:CategoryStatistics* resources are accessed from codes (*skos:Concepts*) via *disco:statisticsCategory*<sup>-</sup>.
  - severity level: ERROR
- **DISCO-C-INVERSE-OBJECT-PROPERTIES-02:** *disco:SummaryStatistics* resources are accessed from *disco:Variables* via *disco:statisticsVariable*<sup>-</sup>.
  - severity level: ERROR
- **DISCO-C-INVERSE-OBJECT-PROPERTIES-03:** *disco:Variables* are accessed from *disco:Questions* via *disco:question*<sup>-</sup>.
  - severity level: ERROR

## 2.7 Symmetric Object Properties\*<sup>B</sup>

A role is symmetric if it is equivalent to its own inverse [5]. An object property symmetry axiom<sup>12</sup> states that the object property expression *OPE* is symmetric - that is, if an individual *x* is connected by *OPE* to an individual *y*, then *y* is also connected by *OPE* to *x*.

<sup>11</sup> *R-56-INVERSE-OBJECT-PROPERTIES*

<sup>12</sup> *R-61-SYMMETRIC-OBJECT-PROPERTIES*

## 2.8 Asymmetric Object Properties\*B

A property is asymmetric<sup>13</sup> if it is disjoint from its own inverse [5]. An object property asymmetry axiom states that the object property *OP* is asymmetric - that is, if an individual *x* is connected by *OP* to an individual *y*, then *y* cannot be connected by *OP* to *x*.

- **DISCO-C-ASYMMETRIC-OBJECT-PROPERTIES-01:** A *disco:Variable* may be based on a *disco:RepresentedVariable*. A *disco:RepresentedVariable*, however, cannot be based on a *disco:Variable*. This is a kind of mistake which may occur as a semantically equivalent object property for the other direction may also be possible (*disco:basisOf*) (*basedOn*  $\sqcap$  *basedOn*<sup>-</sup>  $\sqsubseteq \perp$ ).
  - severity level: ERROR

## 2.9 Reflexive Object Properties\*B

*Reflexive Object Properties*<sup>14</sup> (*reflexive roles*, *global reflexivity* in DL) can be expressed by imposing local reflexivity on the top concept [5].

## 2.10 Irreflexive Object Properties\*B

An object property is irreflexive<sup>15</sup> (*irreflexive role* in DL) if it is never locally reflexive [5]. An object property irreflexivity axiom *IrreflexiveObjectProperty*(*OPE*) states that the object property expression *OPE* is irreflexive - that is, no individual is connected by *OPE* to itself.

- **DISCO-C-IRREFLEXIVE-OBJECT-PROPERTIES-01:** In *Disco*, every object property is irreflexive. No individual is connected by the object property *instrument* to itself ( $\top \sqsubseteq \neg \exists \text{instrument.Self}$ ).
  - severity level: ERROR

## 2.11 Class-Specific Irreflexive Object Properties\*B

A property is *irreflexive* if it is never locally reflexive [5]. An object property irreflexivity axiom states that the object property *OP* is irreflexive - that is, no individual is connected by *OP* to itself. *Class-Specific Irreflexive Object Properties* are object properties which are irreflexive within a given context, e.g. a class.

- **DISCO-C-CLASS-SPECIFIC-IRREFLEXIVE-OBJECT-PROPERTIES-01:** Within the Disco context, *skos:Concepts* cannot be related via the object property *skos:broader* to themselves ( $\text{Concept} \sqsubseteq \neg \exists \text{broader.Self}$ ).
  - severity level: ERROR
- **DISCO-C-CLASS-SPECIFIC-IRREFLEXIVE-OBJECT-PROPERTIES-02:** Within the Disco context, *skos:Concepts* cannot be related via the object property *skos:narrower* to themselves ( $\text{Concept} \sqsubseteq \neg \exists \text{narrower.Self}$ ).
  - severity level: ERROR

<sup>13</sup> R-62-ASYMMETRIC-OBJECT-PROPERTIES

<sup>14</sup> R-59-REFLEXIVE-OBJECT-PROPERTIES

<sup>15</sup> R-60-IRREFLEXIVE-OBJECT-PROPERTIES

## 2.12 Disjoint Properties<sup>B</sup>

A *disjoint properties axiom*<sup>16</sup> states that all of the properties are pairwise disjoint; that is, no individual  $x$  can be connected to an individual/literal  $y$  by these properties.

- **DATA-CUBE-C-DISJOINT-PROPERTIES-01:** All *Data Cube* properties (not having the same domain and range classes) are defined to be pairwise disjoint. The properties *qb:dataSet* and *qb:structure* are disjoint ( $dataSet \sqsubseteq \neg structure$ ).
  - severity level: ERROR
- **DCAT-C-DISJOINT-PROPERTIES-01:** All *DCAT* properties (not having the same domain and range classes) are defined to be pairwise disjoint.
  - severity level: ERROR
- **DISCO-C-DISJOINT-PROPERTIES-01:** All *Disco* properties (not having the same domain and range classes) are defined to be pairwise disjoint. The properties *disco:variable* and *disco:question* are disjoint ( $variable \sqsubseteq \neg question$ ).
  - severity level: ERROR
- **PHDD-C-DISJOINT-PROPERTIES-01:** All *PHDD* properties (not having the same domain and range classes) are defined to be pairwise disjoint. The properties *phdd:isStructuredBy* and *phdd:column* are disjoint ( $isStructuredBy \sqsubseteq \neg column$ ).
  - severity level: ERROR
- **SKOS-C-DISJOINT-PROPERTIES-01:** All *SKOS* properties (not having the same domain and range classes) are defined to be pairwise disjoint.
  - severity level: ERROR
- **SKOS-C-DISJOINT-PROPERTIES-02<sup>17</sup>:** Disjoint Labels Violation: Covers condition S13 from the SKOS reference document stating that “*skos:prefLabel*, *skos:altLabel* and *skos:hiddenLabel* are pairwise disjoint properties”.
  - Implementation: A SPARQL query collects all labels of all concepts, building an in-memory structure. This structure is then checked for disjoint entries.
  - severity level: ERROR
- **XKOS-C-DISJOINT-PROPERTIES-01:** All *XKOS* properties (not having the same domain and range classes) are defined to be pairwise disjoint.
  - severity level: ERROR

<sup>16</sup> *R-9-DISJOINT-PROPERTIES*

<sup>17</sup> Corresponds to qSKOS Quality Issues - SKOS Semi-Formal Consistency Issues - Disjoint Labels Violation

### 2.13 Disjoint Classes<sup>B</sup>

*Disjoint Classes*<sup>18</sup> state that all of the classes are pairwise disjoint; that is, no individual can be at the same time an instance of these disjoint classes.

- **DATA-CUBE-C-DISJOINT-CLASSES-01:** All *Data Cube* classes are defined to be pairwise disjoint.
  - severity level: ERROR
- **DCAT-C-DISJOINT-CLASSES-01:** All *DCAT* classes are defined to be pairwise disjoint.
  - severity level: ERROR
- **DISCO-C-DISJOINT-CLASSES-01:** All *Disco* classes are defined to be pairwise disjoint (e.g. `Study`  $\sqcap$  `Variable`  $\sqsubseteq \perp$ ).
  - severity level: ERROR
- **PHDD-C-DISJOINT-CLASSES-01:** All *PHDD* classes are defined to be pairwise disjoint.
  - severity level: ERROR
- **SKOS-C-DISJOINT-CLASSES-01:** All *SKOS* classes are defined to be pairwise disjoint.
  - severity level: ERROR
- **XKOS-C-DISJOINT-CLASSES-01:** All *XKOS* classes are defined to be pairwise disjoint.
  - severity level: ERROR

### 2.14 Context-Specific Property Groups<sup>C</sup>

The *Context-Specific Property Groups*<sup>19</sup> constraint groups data and object properties within a context (e.g. a class).

### 2.15 Context-Specific Inclusive OR of Properties<sup>C</sup>

*Inclusive or* is a logical connective joining two or more predicates that yields the logical value "true" when at least one of the predicates is true. *Context-Specific Inclusive OR of Properties*<sup>20</sup> constraints specify that individuals are valid if they have at least one property relationship of one or multiple properties stated within a given context. The context can be an application profile, a shape, or a class, i.e., the constraint applies for individuals of this specific class.

<sup>18</sup> *R-7-DISJOINT-CLASSES*

<sup>19</sup> *R-66-PROPERTY-GROUPS*

<sup>20</sup> *R-202-CONTEXT-SPECIFIC-INCLUSIVE-OR-OF-PROPERTIES*



## 2.16 Context-Specific Inclusive OR of Property Groups<sup>C</sup>

At least one property group must match for individuals of a specific context. Context may be a class, a shape, or an application profile.

## 2.17 Recursive Queries<sup>C</sup>

Resource Shapes is a recursive language<sup>21</sup> (the value shape of a Resource Shape is in turn another Resource Shape). There is no way to express that in SPARQL without hand-waving "and then you call the function again here" or "and then you embed this operation here" text. The embedding trick doesn't work in the general case because SPARQL can't express recursive queries, e.g. "test that this Issue is valid and all of the Issues that references, recursively". Most SPARQL engines already have functions that go beyond the official SPARQL 1.1 spec. The cost of that sounds manageable.

## 2.18 Individual Inequality<sup>B</sup>

An *individual inequality axiom*<sup>22</sup> `DifferentIndividuals( a1 ... an )` states that all of the individuals a<sub>i</sub>, 1 ≤ i ≤ n, are different from each other; that is, no individuals a<sub>i</sub> and a<sub>j</sub> with i ≠ j can be derived to be equal. This axiom can be used to axiomatize the unique name assumption — the assumption that all different individual names denote different individuals.

## 2.19 Equivalent Properties<sup>\*B</sup>

An *equivalent object properties axiom*<sup>23</sup> `EquivalentObjectProperties( OPE1 ... OPEn )` states that all of the object property expressions OPE<sub>i</sub>, 1 ≤ i ≤ n, are semantically equivalent to each other. This axiom allows one to use each OPE<sub>i</sub> as a synonym for each OPE<sub>j</sub> — that is, in any expression in the ontology containing such an axiom, OPE<sub>i</sub> can be replaced with OPE<sub>j</sub> without affecting the meaning of the ontology. The axiom `EquivalentObjectProperties( OPE1 OPE2 )` is equivalent to the following two axioms `SubObjectPropertyOf( OPE1 OPE2 )` and `SubObjectPropertyOf( OPE2 OPE1 )`.

An *equivalent data properties axiom*<sup>24</sup> `EquivalentDataProperties( DPE1 ... DPEn )` states that all the data property expressions DPE<sub>i</sub>, 1 ≤ i ≤ n, are semantically equivalent to each other. This axiom allows one to use each DPE<sub>i</sub> as a synonym for each DPE<sub>j</sub> — that is, in any expression in the ontology containing such an axiom, DPE<sub>i</sub> can be replaced with DPE<sub>j</sub> without affecting the meaning of the ontology. The axiom `EquivalentDataProperties( DPE1 DPE2 )` can be seen as a syntactic shortcut for the following axiom `SubDataPropertyOf( DPE1 DPE2 )` and `SubDataPropertyOf( DPE2 DPE1 )`.

<sup>21</sup> R-222-RECURSIVE-QUERIES

<sup>22</sup> R-14-DISJOINT-INDIVIDUALS

<sup>23</sup> R-4-EQUIVALENT-OBJECT-PROPERTIES

<sup>24</sup> R-5-EQUIVALENT-DATA-PROPERTIES

- **DATA-CUBE-C-EQUIVALENT-PROPERTIES-01:** Equivalent properties from different versions of *Data Cube* can be marked as equivalent. As a consequence, the properties can be replaced by each other without affecting the meaning.
- **DCAT-C-EQUIVALENT-PROPERTIES-01:** Equivalent properties from different versions of *DCAT* can be marked as equivalent. As a consequence, the properties can be replaced by each other without affecting the meaning.
- **DISCO-C-EQUIVALENT-PROPERTIES-01:** Equivalent properties from different versions of *Disco* can be marked as equivalent, e.g. *disco:containsVariable* and *disco:variable*. As a consequence, the properties can be replaced by each other without affecting the meaning.
- **PHDD-C-EQUIVALENT-PROPERTIES-01:** Equivalent properties from different versions of *PHDD* can be marked as equivalent. As a consequence, the properties can be replaced by each other without affecting the meaning.
- **SKOS-C-EQUIVALENT-PROPERTIES-01:** Equivalent properties from different versions of *SKOS* can be marked as equivalent. As a consequence, the properties can be replaced by each other without affecting the meaning.
- **XKOS-C-EQUIVALENT-PROPERTIES-01:** Equivalent properties from different versions of *XKOS* can be marked as equivalent. As a consequence, the properties can be replaced by each other without affecting the meaning.

## 2.20 Property Assertions<sup>B</sup>

*Property Assertions*<sup>25</sup> and includes positive property assertions and negative property assertions. A *positive object property assertion* *ObjectPropertyAssertion( OPE a<sub>1</sub> a<sub>2</sub> )* states that the individual a<sub>1</sub> is connected by the object property expression *OPE* to the individual a<sub>2</sub>. A *negative object property assertion* *NegativeObjectPropertyAssertion( OPE a<sub>1</sub> a<sub>2</sub> )* states that the individual a<sub>1</sub> is not connected by the object property expression *OPE* to the individual a<sub>2</sub>. A *positive data property assertion* *DataPropertyAssertion( DPE a lt )* states that the individual a is connected by the data property expression *DPE* to the literal *lt*. A *negative data property assertion* *NegativeDataPropertyAssertion( DPE a lt )* states that the individual a is not connected by the data property expression *DPE* to the literal *lt*.

## 2.21 Data Property Facets<sup>S</sup>

For datatype properties it should be possible to declare frequently needed *facets*<sup>26</sup> to drive user interfaces and validate input against simple conditions, including

<sup>25</sup> *R-96-PROPERTY-ASSERTIONS*

<sup>26</sup> *R-46-CONSTRAINING-FACETS*

min/max value, regular expressions, string length - similar to XSD datatypes. Constraining facets, to restrict datatypes of RDF literals, may be: *xsd:length*, *xsd:minLength*, *xsd:maxLength*, *xsd:pattern*, *xsd:enumeration*, *xsd:whiteSpace*, *xsd:maxInclusive*, *xsd:maxExclusive*, *xsd:minExclusive*, *xsd:minInclusive*, *xsd:totalDigits*, *xsd:fractionDigits*.

- **DISCO-C-DATA-PROPERTY-FACETS-01:** The abstract of a series (*dcterms:abstract*) should have a minimum length (*xsd:minLength*) of some determined minimum length *X*.
  - severity level: WARNING
- **DISCO-C-DATA-PROPERTY-FACETS-02:** The abstract of a study (*dcterms:abstract*) should have a minimum length (*xsd:minLength*) of some determined minimum length *X*.
  - severity level: WARNING

## 2.22 Literal Pattern Matching<sup>S</sup>

There are multiple use cases associated with the requirement to match literals according to given patterns<sup>27</sup>.

- **DISCO-C-LITERAL-PATTERN-MATCHING-01:** Each *disco:Variable* of a given *disco:LogicalDataSet* must have a given prefix for its variable name (*skos:notation*).
  - severity level: INFO

## 2.23 Negative Literal Pattern Matching<sup>S</sup>

Literals of given data properties within given contexts do not have to match given patterns<sup>28</sup>.

- **DISCO-C-NEGATIVE-LITERAL-PATTERN-MATCHING-01:**

## 2.24 Object Property Paths<sup>\*B</sup>

*Object Property Paths*<sup>29</sup> (or *Object Property Chains* and in DL terminology *complex role inclusion axiom* or *role composition*) is the more complex form of sub properties. This axiom states that, if an individual *x* is connected by a sequence of object property expressions *OPE*<sub>1</sub>, ..., *OPE*<sub>*n*</sub> with an individual *y*, then *x* is also connected with *y* by the object property expression *OPE*. Role composition can only appear on the left-hand side of complex role inclusions [5].

<sup>27</sup> R-44-PATTERN-MATCHING-ON-RDF-LITERALS

<sup>28</sup> R-44-PATTERN-MATCHING-ON-RDF-LITERALS

<sup>29</sup> R-55-OBJECT-PROPERTY-PATHS

## 2.25 Intersection\*B

Concept inclusions allow us to state that all mothers are female and that all mothers are parents, but what we really mean is that mothers are exactly the female parents. DLs support such statements by allowing us to form complex concepts such as the *intersection*<sup>30</sup> (also called *conjunction*) which denotes the set of individuals that are both female and parents. A complex concept can be used in axioms in exactly the same way as an atomic concept, e.g., in the equivalence  $\text{Mother} \equiv \text{Female} \sqcap \text{Parent}$ .

## 2.26 Disjunction\*B

A *union class expression*<sup>31</sup> contains all individuals that are instances of at least one class  $C_i$  for  $1 \leq i \leq n$ . A *union data range* contains all tuples of literals that are contained in at least one data range  $DR_i$  for  $1 \leq i \leq n$ . Synonyms of *disjunction* are *union* and *inclusive or*.

- **DISCO-C-DISJUNCTION-01:** Only *disco:Variables* or *disco:Questions* or *disco:RepresentedVariables* can have *disco:concept* relationships to *skos:Concepts*.  
Variable  $\sqcup$  Question  $\sqcup$  RepresentedVariable  $\sqsubseteq \forall$  concept.Concept
  - severity level: ERROR

## 2.27 Negation\*B

A *complement class expression*<sup>32</sup> *ObjectComplementOf( CE )* contains all individuals that are not instances of the class expression *CE*.

## 2.28 Existential Quantifications\*B

An *existential class expression*<sup>33</sup> (*existential restriction* in DL terminology) contains all those individuals that are connected by the property *P* to an individual *x* that is an instance of the class *C* or to literals that are in the data range *DR*.

- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-01:** There must be at least one *disco:universe* relationship from *disco:StudyGroups* to *disco:Universe* (StudyGroup  $\sqsubseteq \exists$  universe.Universe).
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-02:** There must be at least one *disco:universe* relationship from *disco:Studies* to *disco:Universe* (Study  $\sqsubseteq \exists$  universe.Universe).

<sup>30</sup> R-15-CONJUNCTION-OF-CLASS-EXPRESSIONS, R-16-CONJUNCTION-OF-DATA-RANGES

<sup>31</sup> R-17-DISJUNCTION-OF-CLASS-EXPRESSIONS, R-18-DISJUNCTION-OF-DATA-RANGES

<sup>32</sup> R-19-NEGATION-OF-CLASS-EXPRESSIONS, R-20-NEGATION-OF-DATA-RANGES

<sup>33</sup> R-86-EXISTENTIAL-QUANTIFICATION-ON-PROPERTIES

- severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-03:** There may be a *disco:universe* relationship from *disco:RepresentedVariable* to *disco:Universe* (**RepresentedVariable**  $\sqsubseteq \exists$  **universe.Universe**).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-04:** There may be a *disco:universe* relationship from *disco:Variable* to *disco:Universe* (**Variable**  $\sqsubseteq \exists$  **universe.Universe**).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-05:** There may be a *disco:universe* relationship from *disco:Question* to *disco:Universe* (**Question**  $\sqsubseteq \exists$  **universe.Universe**).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-06:** There may be a *disco:universe* relationship from *disco:LogicalDataSet* to *disco:Universe* (**LogicalDataSet**  $\sqsubseteq \exists$  **universe.Universe**).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-07:** There is no relationship (*disco:ddifile*) to a DDI-XML file containing further information about the series (*disco:StudyGroup*) for further analyses.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-08:** There is no relationship (*disco:ddifile*) to a DDI-XML file containing further information about the study (*disco:Study*) for further analyses.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-09:** It is important to know the kind of data (*disco:kindOfData*) collected for a particular series (*disco:StudyGroup*). Survey data, e.g., is much easier accessible than census data. For census data, it is necessary to get in contact with the individual data archive and if data access is granted it may take months to actually get the data.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-10:** It is important to know the kind of data (*disco:kindOfData*) collected for a particular study (*disco:Study*). Survey data, e.g., is much easier accessible than census data. For census data, it is necessary to get in contact with the individual data archive and if data access is granted it may take months to actually get the data.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-11:** Information about the temporal coverage (*dcterms:temporal*) of a series (*disco:StudyGroup*) is of interest for particular queries (e.g. to search for all series of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO

- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-12:** Information about the spatial coverage (*dterms:spatial*) of a series (*disco:StudyGroup*) is of interest for particular queries (e.g. to search for all series of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-13:** Information about the topical coverage (*dterms:subject*) of a series (*disco:StudyGroup*) is of interest for particular queries (e.g. to search for all series of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-14:** Information about the temporal coverage (*dterms:temporal*) of a study (*disco:Study*) is of interest for particular queries (e.g. to search for all studies of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-15:** Information about the spatial coverage (*dterms:spatial*) of a study (*disco:Study*) is of interest for particular queries (e.g. to search for all studies of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-16:** Information about the topical coverage (*dterms:subject*) of a study (*disco:Study*) is of interest for particular queries (e.g. to search for all studies of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-17:** Information about the temporal coverage (*dterms:temporal*) of a data set (*disco:LogicalDataSet*) is of interest for particular queries (e.g. to search for all data sets of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-18:** Information about the spatial coverage (*dterms:spatial*) of a data set (*disco:LogicalDataSet*) is of interest for particular queries (e.g. to search for all data sets of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-19:** Information about the topical coverage (*dterms:subject*) of a data set (*disco:LogicalDataSet*) is of interest for particular queries (e.g. to search for all data sets of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).

- severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-20:** Information about the temporal coverage (*dterms:temporal*) of a data file (*disco:DataFile*) is of interest for particular queries (e.g. to search for all data files of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-21:** Information about the spatial coverage (*dterms:spatial*) of a data file (*disco:DataFile*) is of interest for particular queries (e.g. to search for all data files of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-22:** Information about the topical coverage (*dterms:subject*) of a data file (*disco:DataFile*) is of interest for particular queries (e.g. to search for all data files of a given year (temporal coverage) and for which data is collected in which countries (spatial coverage) about which topics (topical coverage)).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-23:** Information about creators (*dterms:creator*) (persons or organizations) of a series is important when searching for series of the same creators.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-24:** Information about creators (*dterms:creator*) (persons or organizations) of a studies is important when searching for studies of the same creators.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-25:** If summary statistics are collected for studies, detailed further analyses are possible.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-26:** If category statistics are collected for studies, detailed further analyses are possible.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-27:** If a study has no associated data sets, the actual description of the data is missing. Eventually, it is very hard or even impossible to get access to the data.
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-28:** If there is no data file for a given data set, the description of the data set and the containing study is not sufficient.
  - severity level: WARNING
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-29:** The case quantity measures how many cases are collected for a study. High case quantity (*disco:caseQuantity*), stated for data files, is an indicator for high statistical quality of the underlying study. It indicates how comprehensive the study is.

- severity level: WARNING
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-30:** High variable quantity (*disco:variableQuantity*), stated for data files, is an indicator for high statistical quality of the underlying study. It indicates how comprehensive the study is.
  - severity level: WARNING
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-31:** High variable quantity (*disco:variableQuantity*), stated for data sets, is an indicator for high statistical quality of the underlying study. It indicates how comprehensive the study is.
  - severity level: WARNING
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-32:** There is no summary statistics type information (*disco:summaryStatisticsType*) for a summary statistics resource.
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-33:** There is no summary statistics value (*rdf:value*) for a summary statistics resource.
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-34:** There is no relationship to a variable (*disco:statisticsVariable*) for a summary statistics resource.
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-35:** Category statistics resources must be related (*disco:statisticsCategory*) to codes/categories
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-36:** Category statistics resources must have at minimum one value for either frequency, percentage, or cumulative percentage.
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-37:** Codes should be associated with categories (human-readable labels).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-38:** An instrument (*disco:Instrument*) may have a link (*disco:externalDocumentation*) to the questionnaire.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-39:** Questions (*disco:Question*) must have question texts (*disco:questionText*).
  - severity level: ERROR
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-40:** Questions (*disco:Question*) may have response domains (*disco:responseDomain*).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-41:** Questionnaires (*disco:Questionnaire*) may contain (*disco:question*) questions (*disco:Question*).
  - severity level: INFO



- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-42:** Questions (*disco:Question*) may have question numbers (*skos:prefLabel*).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-43:** Variables (*disco:Variable*) may have relationships (*disco:question*) to questions (*disco:Question*), as variables are created out of questions or calculated on the basis of other variables.
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-44:** Data sets (*disco:LogicalDataSet*) may have (*disco:variable*) variables (*disco:Variable*).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-45:** Variables (*disco:Variable*) may have (*disco:concept*) an associated theoretical concept (*skos:Concept*).
  - severity level: INFO
- **DISCO-C-EXISTENTIAL-QUANTIFICATIONS-46:** Each variable (*disco:Variable*) should have (*disco:representation*) a variable representation (*disco:Representation*) which is either an ordered code list (*skos:OrderedCollection*), an unordered code list (*skos:ConceptScheme*) or a union of datatypes (*rdfs:Datatype*).
  - severity level: WARNING
- **DATA-CUBE-C-EXISTENTIAL-QUANTIFICATIONS-01:** Dimensions have range (*IC-4* [3]) - Every dimension declared in a *qb:DataStructureDefinition* must have a declared *rdfs:range*.
  - severity level: ERROR
- **DATA-CUBE-C-EXISTENTIAL-QUANTIFICATIONS-02:** Concept dimensions have code lists (*IC-5* [3]) - Every dimension with range *skos:Concept* must have a *qb:codeList*.
  - severity level: ERROR
- **DATA-CUBE-C-EXISTENTIAL-QUANTIFICATIONS-03:** DSD includes measure (*IC-3* [3]) - Every *qb:DataStructureDefinition* must include (*qb:component*, *qb:componentProperty*) at least one declared measure.
  - severity level: ERROR
- **DATA-CUBE-C-EXISTENTIAL-QUANTIFICATIONS-04 :** Slice Keys must be declared (*IC-7* [3]) - Every *qb:SliceKey* must be associated with (*qb:sliceKey*) a *qb:DataStructureDefinition* (*SliceKey*  $\sqsubseteq \exists$  *sliceKey*<sup>-</sup>.*DataStructureDefinition*).
  - severity level: ERROR

## 2.29 Universal Quantifications\*<sup>B</sup>

A *universal class expression*<sup>34</sup> (*value restriction* in DL) contains all those individuals that are connected by an object property only to individuals that are instances of a particular class.

- **DATA-CUBE-C-UNIVERSAL-QUANTIFICATIONS-01:** *Universal quantifications* are defined for each *Data Cube* object and data property.

<sup>34</sup> R-91-UNIVERSAL-QUANTIFICATION-ON-PROPERTIES

- Severity level: ERROR
- **DCAT-C-UNIVERSAL-QUANTIFICATIONS-01:** *Universal quantifications* are defined for each *DCAT* object and data property. Only *dcat:Catalogs* can have *dcat:dataset* relationships to *dcat:Datasets* ( $\text{Catalog} \sqsubseteq \forall \text{ dataset.Dataset}$ ).
  - Severity level: ERROR
- **DISCO-C-UNIVERSAL-QUANTIFICATIONS-01:** *Universal quantifications* are defined for each *Disco* object and data property. Only *disco:LogicalDataSets* can have *disco:aggregation* relationships to *qb:DataSets* ( $\text{LogicalDataSet} \sqsubseteq \forall \text{ aggregation.DataSet}$ ).
  - Severity level: ERROR
- **PHDD-C-UNIVERSAL-QUANTIFICATIONS-01:** *Universal quantifications* are defined for each *PHDD* object and data property.
  - Severity level: ERROR
- **SKOS-C-UNIVERSAL-QUANTIFICATIONS-01:** *Universal quantifications* are defined for each *SKOS* object and data property.
  - Severity level: ERROR
- **XKOS-C-UNIVERSAL-QUANTIFICATIONS-01:** *Universal quantifications* are defined for each *XKOS* object and data property.
  - Severity level: ERROR

### 2.30 Minimum Unqualified Cardinality Restrictions\*B

A *minimum cardinality restriction*<sup>35</sup> contains all those individuals that are connected by a property to at least  $n$  different individuals/literals that are instances of a particular class or data range. If the class is missing, it is taken to be *owl:Thing*. If the data range is missing, it is taken to be *rdfs:Literal*.  $\leq nR$ .  $\top$  is the minimum unqualified cardinality restriction where  $n \in \mathbb{N}$  (written  $\leq nR$  in short). For unqualified cardinality restrictions, classes respective data ranges are not stated.

### 2.31 Minimum Qualified Cardinality Restrictions\*B

A *minimum cardinality restriction*<sup>36</sup> contains all those individuals that are connected by a property to at least  $n$  different individuals/literals that are instances of a particular class or data range. If the class is missing, it is taken to be *owl:Thing*. If the data range is missing, it is taken to be *rdfs:Literal*.  $\geq nR$ .  $C$  is a minimum qualified cardinality restriction where  $n \in \mathbb{N}$ .

<sup>35</sup> R-81-MINIMUM-UNQUALIFIED-CARDINALITY-ON-PROPERTIES, R-211-CARDINALITY-CONSTRAINTS

<sup>36</sup> R-75-MINIMUM-QUALIFIED-CARDINALITY-ON-PROPERTIES, R-211-CARDINALITY-CONSTRAINTS

- **DATA-CUBE-C-MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** *Minimum Qualified Cardinality Restrictions* constraints are defined for each *Data Cube* object and data property.
  - Severity level: ERROR
- **DATA-CUBE-C-MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-02:** Unique data set (*IC-1* [3]) - Every *qb:Observation* has (*qb:dataSet*) exactly one associated *qb:DataSet* (*Observation*  $\sqsubseteq \geq 1$  *dataSet.DataSet*  $\sqcap \leq 1$  *dataSet.DataSet*).
  - Severity level: ERROR
- **DCAT-C-MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** *Minimum Qualified Cardinality Restrictions* constraints are defined for each *DCAT* object and data property.
  - Severity level: ERROR
- **DISCO-C-MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** *Minimum Qualified Cardinality Restrictions* constraints are defined for each *Disco* object and data property. A *disco:Questionnaire*, e.g., has at least one *disco:question* relationship to *disco:Questions* (*Questionnaire*  $\sqsubseteq \geq 1$  *question.Question*).
  - Severity level: ERROR
- **PHDD-C-MINIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** *Minimum Qualified Cardinality Restrictions* constraints are defined for each *PHDD* object and data property.
  - Severity level: ERROR

### 2.32 Maximum Unqualified Cardinality Restrictions\*<sup>B</sup>

A *maximum cardinality restriction* contains all those individuals that are connected by a property to at most *n* different individuals/literals that are instances of a particular class or data range. If the class is missing, it is taken to be *owl:Thing*. If the data range is not present, it is taken to be *rdfs:Literal*. Unqualified means that the class respective the data range is not stated.  $\geq nR.T$  is a *maximum unqualified cardinality restriction*<sup>37</sup> where  $n \in \mathbb{N}$  (written  $\geq nR$  in short).

### 2.33 Maximum Qualified Cardinality Restrictions\*<sup>B</sup>

A *maximum cardinality restriction* contains all those individuals that are connected by a property to at most *n* different individuals/literals that are instances of a particular class or data range. If the class is missing, it is taken to be *owl:Thing*. If the data range is not present, it is taken to be *rdfs:Literal*. Qualified means that the class respective the data range is stated.  $\leq nR.C$  is a *maximum qualified cardinality restriction*<sup>38</sup> where  $n \in \mathbb{N}$ .

<sup>37</sup> *R-82-MAXIMUM-UNQUALIFIED-CARDINALITY-ON-PROPERTIES, R-211-CARDINALITY-CONSTRAINTS*

<sup>38</sup> *R-76-MAXIMUM-QUALIFIED-CARDINALITY-ON-PROPERTIES, R-211-CARDINALITY-CONSTRAINTS*

- **DISCO-C-MAXIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** A *disco:Variable* has at most one *disco:concept* relationship to a theoretical concept (*skos:Concept*) ( $\text{Variable} \sqsubseteq \leq 1 \text{ concept.Concept}$ ).
  - Severity level: ERROR
- **DATA-CUBE-C-MAXIMUM-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** Unique data set (*IC-1* [3]) - Every *qb:Observation* has (*qb:dataSet*) exactly one associated *qb:DataSet* ( $\text{Observation} \sqsubseteq \geq 1 \text{ dataSet.DataSet} \sqcap \leq 1 \text{ dataSet.DataSet}$ ).
  - Severity level: ERROR

### 2.34 Exact Unqualified Cardinality Restrictions\*B

An *exact cardinality restriction*<sup>39</sup> contains all those individuals that are connected by a property to exactly *n* different individuals that are instances of a particular class or data range. If the class is missing, it is taken to be *owl:Thing*. If the data range is not present, it is taken to be *rdfs:Literal*. Unqualified means that the class respective data range is not stated.  $\geq nR.\top \sqcap \leq nR.\top$  is an exact unqualified cardinality restriction where  $n \in \mathbb{N}$ .

- **DATA-CUBE-C-EXACT-UNQUALIFIED-CARDINALITY-RESTRICTIONS-01:** Unique slice structure (*IC-9* [3]) - Each *qb:Slice* must have exactly one associated *qb:sliceStructure*.
  - Severity level: ERROR

### 2.35 Exact Qualified Cardinality Restrictions\*B

An *exact cardinality restriction*<sup>40</sup> contains all those individuals that are connected by a property to exactly *n* different individuals that are instances of a particular class or data range. If the class is missing, it is taken to be *owl:Thing*. If the data range is not present, it is taken to be *rdfs:Literal*.  $\geq nR.C \sqcap \leq nR.C$  is an exact qualified cardinality restriction where  $n \in \mathbb{N}$ .

- **DISCO-C-EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-01:** A *disco:Question* has exactly 1 *disco:universe* relationship to *disco:Universe* ( $\text{Question} \sqsubseteq \geq 1 \text{ universe.Universe} \sqcap \leq 1 \text{ universe.Universe}$ ).
  - Severity level: ERROR
- **DATA-CUBE-C-EXACT-QUALIFIED-CARDINALITY-RESTRICTIONS-02:** Unique DSD (*IC-2* [3]) - Every *qb:DataSet* has (*qb:structure*) exactly one associated *qb:DataStructureDefinition* ( $\text{DataSet} \sqsubseteq \geq 1 \text{ structure.DataStructureDefinition} \sqcap \leq 1 \text{ structure.DataStructureDefinition}$ ).
  - Severity level: ERROR

<sup>39</sup> *R-80-EXACT-UNQUALIFIED-CARDINALITY-ON-PROPERTIES, CARDINALITY-CONSTRAINTS* R-211-

<sup>40</sup> *R-74-EXACT-QUALIFIED-CARDINALITY-ON-PROPERTIES, CARDINALITY-CONSTRAINTS* R-211-

### 2.36 Transitive Object Properties\*B

*Transitivity* is a special form of *complex role inclusion*. An *object property transitivity axiom*<sup>41</sup> states that the object property is transitive — that is, if an individual  $x$  is connected by the object property to an individual  $y$  that is connected by the object property to an individual  $z$ , then  $x$  is also connected by the object property to  $z$ .

### 2.37 Context-Specific Exclusive OR of Properties<sup>C</sup>

*Exclusive or* is a logical operation that outputs true whenever both inputs differ (one is true, the other is false). Only one of multiple properties within some context (e.g. a class, a shape, or an application profile) leads to valid data<sup>42</sup>. This constraint is generally expressed in DL as follows:  $C \sqsubseteq (\neg A \sqcap B) \sqcup (A \sqcap \neg B)$ .

### 2.38 Context-Specific Exclusive OR of Property Groups<sup>C</sup>

*Exclusive or* is a logical operation that outputs true whenever both inputs differ (one is true, the other is false). Only one of multiple property groups leads to valid data<sup>43</sup>.

- **DISCO-C-CONTEXT-SPECIFIC-EXCLUSIVE-OR-OF-PROPERTY-GROUPS-01:** Within the context of *Disco*, *skos:Concepts* can have either *skos:definition* (when interpreted as theoretical concepts) or *skos:notation* and *skos:prefLabel* properties (when interpreted as codes and categories),

$$\text{Concept} \sqsubseteq (\neg D \sqcap C) \sqcup (D \sqcap \neg C)$$

$$D \equiv A \sqcap B$$

but not both.  $A \sqsubseteq \geq 1 \text{ notation.string} \sqcap \leq 1 \text{ notation.string}$

$$B \sqsubseteq \geq 1 \text{ prefLabel.string} \sqcap \leq 1 \text{ prefLabel.string}$$

$$C \sqsubseteq \geq 1 \text{ definition.string} \sqcap \leq 1 \text{ definition.string}$$

- severity level: INFO

### 2.39 Allowed Values<sup>B</sup>

It is a common requirement to narrow down the value space of a property by an exhaustive enumeration of the valid values (both literals or resources). This is often rendered in drop down boxes or radio buttons in user interfaces. *Allowed values*<sup>44</sup> for properties can be IRIs, IRIs (matching one or multiple patterns), (any) literals, literals of a list of allowed literals (e.g. 'red' 'blue' 'green'), typed literals of one or multiple type(s) (e.g. *xsd:string*).

<sup>41</sup> R-63-TRANSITIVE-OBJECT-PROPERTIES

<sup>42</sup> R-11-CONTEXT-SPECIFIC-EXCLUSIVE-OR-OF-PROPERTIES

<sup>43</sup> R-13-DISJOINT-GROUP-OF-PROPERTIES-CLASS-SPECIFIC

<sup>44</sup> R-30-ALLOWED-VALUES-FOR-RDF-OBJECTS and R-37-ALLOWED-VALUES-FOR-RDF-LITERALS

- **DISCO-C-ALLOWED-VALUES-01.** *disco:CategoryStatistics* can only have *disco:computationBase* relationships to the values *valid* and *invalid* of the datatype *rdf:langString* ( $\text{CategoryStatistics} \equiv \forall \text{computationBase}.\{\text{valid}, \text{invalid}\} \sqcap \text{langString}$ ).
  - severity level: ERROR

## 2.40 Not Allowed Values<sup>B</sup>

A matching triple has any literal / object except those explicitly excluded<sup>45</sup>.

## 2.41 Literal Ranges<sup>S</sup>

*P1* is a data property (of an instance of class *C1*) and its literal value must be between the range of  $[V_{min}, V_{max}]$ <sup>46</sup>.

- **DISCO-C-LITERAL-RANGES-01:** *disco:percentage* (domain: *disco:CategoryStatistics*) literals must be of the datatype *xsd:double* whose range should be restricted to be between 0 and 100.
  - severity level: ERROR
- **DISCO-C-LITERAL-RANGES-02:** *disco:cumulativePercentage* (domain: *disco:CategoryStatistics*) literals must be of the datatype *xsd:double* whose range should be restricted to be between 0 and 100.
  - severity level: ERROR

## 2.42 Negative Literal Ranges<sup>S</sup>

*P1* is a data property (of an instance of class *C1*) and its literal value must not be between the range of  $[V_{min}, V_{max}]$ <sup>47</sup>.

## 2.43 Required Properties<sup>B</sup>

Properties may be required<sup>48</sup>.

## 2.44 Optional Properties<sup>B</sup>

Properties may be optional<sup>49</sup>.

<sup>45</sup> *R-33-NEGATIVE-OBJECT-CONSTRAINTS*, *R-200-NEGATIVE-LITERAL-CONSTRAINTS*

<sup>46</sup> *R-45-RANGES-OF-RDF-LITERAL-VALUES*

<sup>47</sup> *R-142-NEGATIVE-RANGES-OF-RDF-LITERAL-VALUES*

<sup>48</sup> *R-68-REQUIRED-PROPERTIES*

<sup>49</sup> *R-69-OPTIONAL-PROPERTIES*

## 2.45 Repeatable Properties<sup>B</sup>

Properties may be repeatable<sup>50</sup>.

## 2.46 Negative Property Constraints<sup>S</sup>

Instances of a specific class must not have some object property<sup>51</sup>.

## 2.47 Individual Equality\*<sup>B</sup>

*Individual equality*<sup>52</sup> states that two different names are known to refer to the same individual [5].

## 2.48 Functional Properties\*<sup>B</sup>

An *object property functionality axiom*<sup>53</sup> *FunctionalObjectProperty( OPE )* states that the object property expression *OPE* is functional — that is, for each individual *x*, there can be at most one distinct individual *y* such that *x* is connected by *OPE* to *y*. Each such axiom can be seen as a syntactic shortcut for the following axiom: *SubClassOf( owl:Thing ObjectMaxCardinality( 1 OPE ) )*.

## 2.49 Inverse-Functional Properties\*<sup>B</sup>

An *object property inverse functionality axiom*<sup>54</sup> *InverseFunctionalObjectProperty( OPE )* states that the object property expression *OPE* is inverse-functional - that is, for each individual *x*, there can be at most one individual *y* such that *y* is connected by *OPE* with *x*. Each such axiom can be seen as a syntactic shortcut for the following axiom: *SubClassOf( owl:Thing ObjectMaxCardinality( 1 ObjectInverseOf( OPE ) ) )*.

- **DISCO-C-INVERSE-FUNCTIONAL-PROPERTIES-01:** For each *rdfs:Resource x*, there can be at most one distinct *rdfs:Resource y* such that *y* is connected by *adms:identifier* to *x* (`funcnt identifier`).
  - severity level: ERROR
- **DISCO-C-INVERSE-FUNCTIONAL-PROPERTIES-02:** Keys are even more general than inverse-functional properties, as a key can be a data, an object property, or a chain of properties [7]. 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 a special *keyfor* construct (`identifier keyfor Resource`) [6]. OWL 2 *hasKey* implements *keyfor* and thus can be used to identify resources uniquely, to merge resources with identical key property values, and to recognize constraint violations.
  - severity level: ERROR

<sup>50</sup> R-70-REPEATABLE-PROPERTIES

<sup>51</sup> R-52-NEGATIVE-OBJECT-PROPERTY-CONSTRAINTS, R-53-NEGATIVE-DATA-PROPERTY-CONSTRAINTS

<sup>52</sup> R-6-EQUIVALENT-INDIVIDUALS

<sup>53</sup> R-57-FUNCTIONAL-OBJECT-PROPERTIES

<sup>54</sup> R-58-INVERSE-FUNCTIONAL-OBJECT-PROPERTIES

## 2.50 Value Restrictions\*B

*Individual Value Restrictions*<sup>55</sup>: A has-value class expression *ObjectHasValue( OPE a )* consists of an object property expression *OPE* and an individual *a*, and it contains all those individuals that are connected by *OPE* to *a*. Each such class expression can be seen as a syntactic shortcut for the class expression *ObjectSomeValuesFrom( OPE ObjectOneOf( a ) )*. *Literal Value Restrictions*: A has-value class expression *DataHasValue( DPE lt )* consists of a data property expression *DPE* and a literal *lt*, and it contains all those individuals that are connected by *DPE* to *lt*. Each such class expression can be seen as a syntactic shortcut for the class expression *DataSomeValuesFrom( DPE DataOneOf( lt ) )*.

## 2.51 Self Restrictions\*B

A *self-restriction* *ObjectHasSelf( OPE )* consists of an object property expression *OPE*, and it contains all those individuals that are connected by *OPE* to themselves.

## 2.52 Primary Key Properties\*B

The *Primary Key Properties*<sup>56</sup> constraint is often useful to declare a given (datatype) property as the "primary key" of a class, so that a system can enforce uniqueness and also automatically build URIs from user input and data imported from relational databases or spreadsheets.. Starfleet officers, e.g., are uniquely identified by their command authorization code (e.g. to activate and cancel auto-destruct sequences). It means that the property *commandAuthorizationCode* is inverse functional - mapped to DL as follows: `(func commandAuthorizationCode)`. Keys, however, are even more general, i.e., a generalization of inverse functional properties [7]. A key can be a datatype property, an object property, or a chain of properties. 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 a special *keyfor* construct[6]. This leads to the following DL mapping (only one simple property constraint): `commandAuthorizationCode keyfor StarfleetOfficer`

– see *inverse-functional properties*

## 2.53 Class-Specific Property Range\*B

*Class-Specific Property Range*<sup>57</sup> restricts the range of object and data properties for individuals within a specific context (e.g. class, shape, application profile). The values of each member property of a class may be limited by their value type, such as *xsd:string* or *foaf:Person*.

<sup>55</sup> R-88-VALUE-RESTRICTIONS

<sup>56</sup> R-226-PRIMARY-KEY-PROPERTIES

<sup>57</sup> R-29-CLASS-SPECIFIC-RANGE-OF-RDF-OBJECTS, R-36-CLASS-SPECIFIC-RANGE-OF-RDF-LITERALS



- **DISCO-C-CLASS-SPECIFIC-PROPERTY-RANGE-01**: Only *disco:Questions* can have *disco:questionText* relationships to literals of the datatype *rdf:langString* ( $\neg \text{Question} \sqsubseteq \neg \exists \text{questionText}.\text{langString}$ ).
  - severity level: ERROR

## 2.54 Class-Specific Reflexive Object Properties\*B

Using DL terminology *Class-Specific Reflexive Object Properties* is called local reflexivity - a set of individuals (of a specific class) that are related to themselves via a given role [5].

## 2.55 Membership in Controlled Vocabularies<sup>S</sup>

Resources can only be members of listed controlled vocabularies<sup>58</sup>.

- **DISCO-C-MEMBERSHIP-IN-CONTROLLED-VOCABULARIES-01**: *disco:SummaryStatistics* can only have *disco:summaryStatisticType* relationships to *skos:Concepts* which must be members of the controlled vocabulary *ddicv:SummaryStatisticType* which is a *skos:ConceptScheme*.
 
$$\text{SummaryStatistics} \sqsubseteq \forall \text{summaryStatisticType}.A$$

$$A \equiv \text{Concept} \sqcap \forall \text{inScheme}.B$$

$$B \equiv \text{ConceptScheme} \sqcap \{\text{SummaryStatisticType}\}$$
  - severity level: ERROR
- **DATA-CUBE-C-MEMBERSHIP-IN-CONTROLLED-VOCABULARIES-01**: Codes from code list (*IC-19* [3]) - If a dimension property has a *qb:codeList*, then the value of the dimension property on every *qb:Observation* must be in the code list.
  - severity level: ERROR

## 2.56 IRI Pattern Matching<sup>S</sup>

IRI pattern matching applied on subjects, properties, and objects<sup>59</sup>.

- **DISCO-C-IRI-PATTERN-MATCHING-01**: *disco:Study* resources must match a given IRI pattern.
  - severity level: INFO

<sup>58</sup> *R-32-MEMBERSHIP-OF-RDF-OBJECTS-IN-CONTROLLED-VOCABULARIES*, *R-39-MEMBERSHIP-OF-RDF-LITERALS-IN-CONTROLLED-VOCABULARIES*

<sup>59</sup> *R-21-IRI-PATTERN-MATCHING-ON-RDF-SUBJECTS*, *R-22-IRI-PATTERN-MATCHING-ON-RDF-OBJECTS*, *R-23-IRI-PATTERN-MATCHING-ON-RDF-PROPERTIES*

## 2.57 Literal Value Comparison<sup>S</sup>

Depending on the property semantics, there are cases where two different literal values must have a specific ordering with respect to an operator. *P1* and *P2* are the datatype properties we need to compare and *OP* is the comparison operator (<, <=, >, >=, =, !=)<sup>60</sup>. The *COMP Pattern*, one of the Data Quality Test Patterns, can be used to validate the *Literal Value Comparison* constraint [4]:

```
1 SELECT ?s WHERE {  
2   ?s %%P1%% ?v1 .  
3   ?s %%P2%% ?v2 .  
4   FILTER ( ?v1 %%OP%% ?v2 ) }
```

- **DISCO-C-LITERAL-VALUE-COMPARISON-01**: *disco:startDates* must be before (<) *disco:endDates*. To validate this constraint we bind the variables as follows (P1: *disco:startDate*, P2: *disco:endDate*, OP: <).
  - severity level: ERROR

## 2.58 Ordering<sup>C</sup>

With this constraint objects of object properties can be ordered as well as literals of data properties<sup>61</sup>.

In DDI, variables, questions, and codes/categories are typically organized in a particular order. For obtaining this order, *skos:OrderedCollection* resources are used.

- **DISCO-C-ORDERING-01**: If *disco:Variables* of a given *disco:LogicalDataSet* should be ordered, a collection of variables must be present in the data and connected with the data set. The collection of variables is of the type *skos:OrderedCollection* containing multiple variables (each represented as *skos:Concept*) in a *skos:memberList*.
  - severity level: INFO
- **DISCO-C-ORDERING-02**: If *disco:Questions* of a given *disco:Questionnaire* should be ordered, a collection of questions must be present in the data and connected with the questionnaire. The collection of questions is of the type *skos:OrderedCollection* containing multiple questions (each represented as *skos:Concept*) in a *skos:memberList*.
  - severity level: INFO
- **DISCO-C-ORDERING-03**: If codes/categories (*skos:Concepts*) of a given *disco:Representation* of a given *disco:Variable* should be ordered, the variable representation should also be of the type *skos:OrderedCollection* containing multiple codes/categories (each represented as *skos:Concept*) in a *skos:memberList*.
  - severity level: INFO

<sup>60</sup> R-43-LITERAL-VALUE-COMPARISON

<sup>61</sup> R-121-SPECIFY-ORDER-OF-RDF-RESOURCES, R-217-DEFINE-ORDER-FORMS/DISPLAY

## 2.59 Validation Levels<sup>S</sup>

Different levels of severity (priority)<sup>62</sup> should be assigned to constraints. Possible validation levels could be: informational, warning, error, fail, should, recommended, must, may, optional, closed (only this) constraints, open (at least this) constraint.

For *Disco* each constraint should be assigned to exactly one *validation level*.

## 2.60 String Operations<sup>S</sup>

Many different *string operations*<sup>63</sup> are possible. Some constraints require building new strings out of other strings. Calculating the string length would also be another constraint of this type.

- **DISCO-C-STRING-OPERATIONS-01**: The title of a study (*dcterms:title*) (e.g. 'EU-SILC 2005') may be calculated out of the title of the containing series (*dcterms:title*) (e.g. 'EU-SILC') and the human-readable label of the study (*rdfs:label*) (e.g. '2005').
  - severity level: INFO

## 2.61 Context-Specific Valid Classes<sup>B</sup>

What types are valid in a specific context?<sup>64</sup> Context can be an input stream, a data creation function, or an API.

- **DATA-CUBE-C-CONTEXT-SPECIFIC-VALID-CLASSES-01**: For future versions of *Data Cube*, out-dated classes can be marked as deprecated.
  - severity level: INFO
- **DCAT-C-CONTEXT-SPECIFIC-VALID-CLASSES-01**: For future versions of *DCAT*, out-dated classes can be marked as deprecated.
  - severity level: INFO
- **DISCO-C-CONTEXT-SPECIFIC-VALID-CLASSES-01**: For future versions of *Disco*, out-dated classes can be marked as deprecated.
  - severity level: INFO
- **PHDD-C-CONTEXT-SPECIFIC-VALID-CLASSES-01**: For future versions of *PHDD*, out-dated classes can be marked as deprecated.
  - severity level: INFO
- **SKOS-C-CONTEXT-SPECIFIC-VALID-CLASSES-01**: For future versions of *SKOS*, out-dated classes can be marked as deprecated.
  - severity level: INFO
- **XKOS-C-CONTEXT-SPECIFIC-VALID-CLASSES-01**: For future versions of *XKOS*, out-dated classes can be marked as deprecated.
  - severity level: INFO

<sup>62</sup> *R-205-VARYING-LEVELS-OF-ERROR*, *R-135-CONSTRAINT-LEVELS*, *R-158-SEVERITY-LEVELS-OF-CONSTRAINT-VIOLATIONS*, *R-193-MULTIPLE-CONSTRAINT-VALIDATION-EXECUTION-LEVELS*

<sup>63</sup> *R-194-PROVIDE-STRING-FUNCTIONS-FOR-RDF-LITERALS*

<sup>64</sup> *R-209-VALID-CLASSES*

## 2.62 Context-Specific Valid Properties<sup>B</sup>

What properties can be used within this context?<sup>65</sup> Context can be an data receipt function, data creation function, or API.

- ***DATA-CUBE-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01***: For future versions of *Data Cube*, out-dated properties can be marked as deprecated.
  - severity level: INFO
- ***DCAT-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01***: For future versions of *DCAT*, out-dated properties can be marked as deprecated.
  - severity level: INFO
- ***DISCO-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01***: For future versions of *Disco*, out-dated properties can be marked as deprecated.
  - severity level: INFO
- ***PHDD-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01***: For future versions of *PHDD*, out-dated properties can be marked as deprecated.
  - severity level: INFO
- ***SKOS-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01***: For future versions of *SKOS*, out-dated properties can be marked as deprecated.
  - severity level: INFO
- ***XKOS-C-CONTEXT-SPECIFIC-VALID-PROPERTIES-01***: For future versions of *XKOS*, out-dated properties can be marked as deprecated.
  - severity level: INFO

## 2.63 Default Values<sup>\*S</sup>

*Default values*<sup>66</sup> for objects and literals are inferred automatically. It should be possible to declare the default value for a given property, e.g. so that input forms can be pre-populated and to insert a required property that is missing in a web service call.

- ***DISCO-C-DEFAULT-VALUES-01***: The value 'true' for the property *disco:isPublic* (*xsd:boolean*) indicates that the data set (*disco:LogicalDataSet*) can be accessed (usually downloaded) by anyone. Per default, access to data sets should be restricted ('false').
  - severity level: INFO

<sup>65</sup> *R-210-VALID-PROPERTIES*

<sup>66</sup> *R-31-DEFAULT-VALUES-OF-RDF-OBJECTS*,     *R-38-DEFAULT-VALUES-OF-RDF-LITERALS*

## 2.64 Mathematical Operations<sup>C</sup>

Examples for *Mathematical Operations*<sup>67</sup> are the addition of two dates, the addition of days to a start date, and statistical computations (e.g. average, mean, sum).

- **DISCO-C-MATHEMATICAL-OPERATIONS-01:** The sum of *disco:percentage* (datatype: *xsd:double*) values of all codes (represented as *skos:Concepts*) of a code list (*skos:ConceptScheme* or *skos:OrderedCollection*), serving as representation of a particular *disco:Variable*, must exactly be 100.
  - severity level: ERROR
- **DISCO-C-MATHEMATICAL-OPERATIONS-02:** For a given variable, the sum of the frequencies of all codes of the variable's code list has to be equal to the variable's total number of cases (summary statistics of the type 'number of cases').
  - severity level: ERROR
- **DISCO-C-MATHEMATICAL-OPERATIONS-03:** For a given variable, the sum of 'valid cases' and 'invalid cases' has to be equal to the total 'number of cases'.
  - severity level: ERROR
- **DISCO-C-MATHEMATICAL-OPERATIONS-04:** For a given variable, the total 'number of cases' value for the country 'All' must be equal to the sum of the total 'number of cases' value for each country.
  - severity level: ERROR
- **DISCO-C-MATHEMATICAL-OPERATIONS-05:** Minimum values do not have to be greater than maximum values (*disco:SummaryStatistics*).
  - severity level: ERROR

## 2.65 Language Tag Matching<sup>S</sup>

For particular data properties, values must be stated for predefined languages<sup>68</sup>.

- **DISCO-C-LANGUAGE-TAG-MATCHING-01:** There must be an English variable name (*skos:notation*) for each *disco:Variable* within *disco:LogicalDataSets*.
  - severity level: INFO

## 2.66 Language Tag Cardinality<sup>S</sup>

For particular data properties, values of predefined languages must be stated for determined number of times<sup>69</sup>.

- **DISCO-C-LANGUAGE-TAG-CARDINALITY-01:** There must be at least one English *disco:questionText* for each *disco:Question* within *disco:LogicalDataSets*.

<sup>67</sup> *R-42-MATHEMATICAL-OPERATIONS*, *R-41-STATISTICAL-COMPUTATIONS*

<sup>68</sup> *R-47-LANGUAGE-TAG-MATCHING*

<sup>69</sup> *R-49-RDF-LITERALS-HAVING-AT-MOST-ONE-LANGUAGE-TAG*, *R-48-MISSING-LANGUAGE-TAGS*

- severity level: INFO
- **DISCO-C-LANGUAGE-TAG-CARDINALITY-02:** There should be at most one English literal value for variable names (*skos:notation*, domain: *disco:Variable*).
  - severity level: INFO
- **DISCO-C-LANGUAGE-TAG-CARDINALITY-03:** For each question (*disco:Question*), there must be at least one question text (*disco:questionText*) associated with a language tag of an arbitrary language or with an English language tag.
  - severity level: INFO
- **SKOS-C-LANGUAGE-TAG-CARDINALITY-01<sup>70</sup>:** Omitted or Invalid Language Tags: Some controlled vocabularies contain literals in natural language, but without information what language has actually been used. Language tags might also not conform to language standards, such as RFC 3066.
  - Implementation: Iteration over all triples in the vocabulary that have a predicate which is a (subclass of) *rdfs:label* or *skos:note*.
  - Severity level: WARNING
- **SKOS-C-LANGUAGE-TAG-CARDINALITY-02<sup>71</sup>:** Incomplete Language Coverage: Some concepts in a thesaurus are labeled in only one language, some in multiple languages. It may be desirable to have each concept labeled in each of the languages that also are used on the other concepts. This is not always possible, but incompleteness of language coverage for some concepts can indicate shortcomings of the vocabulary.
  - Severity level: INFO
- **SKOS-C-LANGUAGE-TAG-CARDINALITY-03<sup>72</sup>:** No Common Language: Checks if all concepts have at least one common language, i.e. they have assigned at least one literal in the same language.
  - Severity level: INFO
- **SKOS-C-LANGUAGE-TAG-CARDINALITY-04<sup>73</sup>:** Inconsistent Preferred Labels: According to the SKOS reference document, "A resource has no more than one value of *skos:prefLabel* per language tag".
  - Implementation: A SPARQL query is used to find concepts with at least two *prefLabels*. In a second step, the language tags of these *prefLabels* are analyzed and an ambiguity is detected if they are equal.
  - Severity level: INFO

<sup>70</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Omitted or Invalid Language Tags

<sup>71</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Incomplete Language Coverage

<sup>72</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - No Common Language

<sup>73</sup> Corresponds to qSKOS Quality Issues - SKOS Semi-Formal Consistency Issues - Inconsistent Preferred Labels

## 2.67 Whitespace Handling<sup>S</sup>

Avoid whitespaces in literals neither leading nor trailing white spaces<sup>74</sup>.

- **DISCO-C-WHITESPACE-HANDLING-01**: Delete whitespaces of series and study abstracts (*dcterms:abstract*; domain: *disco:StudyGroup*, *disco:Study*) automatically.
  - severity level: INFO

## 2.68 HTML Handling<sup>S</sup>

Check if all HTML tags, included in literals (of specific data properties within the context of specific classes)<sup>75</sup>, are closed properly.

- **DISCO-C-HTML-HANDLING-01**: Check if all HTML tags, included in literals of all *Disco* data properties, are closed properly.
  - severity level: INFO
- **DISCO-C-HTML-HANDLING-02**: Check if all HTML tags, included in literals of all data properties whose domains are *Disco* classes, are closed properly.
  - severity level: INFO

## 2.69 Conditional Properties<sup>C</sup>

If specific properties exist, then specific other properties must also be present<sup>76</sup>.

- **DISCO-C-CONDITIONAL-PROPERTIES-01**: If a *skos:Concept* represents a code (having a *skos:notation* property) and a category (having a *skos:prefLabel* property), then the property *disco:isValid* has to be stated indicating if the code is valid ('true') or missing ('false').
  - severity level: ERROR
- **DISCO-C-CONDITIONAL-PROPERTIES-02**: In order to get an overview over a series or a study either an abstract, a title, an alternative title, or links to external descriptions should be stated. If the abstract (*dcterms:abstract*) of a series (*disco:StudyGroup*) and an external description of the series (*disco:ddifile*) is missing, a series title (*dcterms:title*) or an alternative series title (*dcterms:alternative*) has to be stated.
  - severity level: WARNING
- **DISCO-C-CONDITIONAL-PROPERTIES-03**: In order to get an overview over a series or a study either an abstract, a title, an alternative title, or links to external descriptions should be stated. If the abstract (*dcterms:abstract*) of a study (*disco:Study*) and an external description of the study (*disco:ddifile*) is missing, a study title (*dcterms:title*) or an alternative study title (*dcterms:alternative*) has to be stated.

<sup>74</sup> R-50-WHITESPACE-HANDLING-OF-RDF-LITERALS

<sup>75</sup> R-51-HTML-HANDLING-OF-RDF-LITERALS

<sup>76</sup> R-71-CONDITIONAL-PROPERTIES

- severity level: WARNING
- ***DISCO-C-CONDITIONAL-PROPERTIES-04***: If the abstract (*dcterms:abstract*) of a series (*disco:StudyGroup*), an external description of the series (*disco:ddifile*), a series title (*dcterms:title*), and an alternative series title (*dcterms:alternative*) is missing, an error message should be shown.
  - severity level: ERROR
- ***DISCO-C-CONDITIONAL-PROPERTIES-05***: If the abstract (*dcterms:abstract*) of a study (*disco:Study*), an external description of the study (*disco:ddifile*), a study title (*dcterms:title*), and an alternative study title (*dcterms:alternative*) is missing, an error message should be shown.
  - severity level: ERROR
- ***DISCO-C-CONDITIONAL-PROPERTIES-06***: If a category statistics resource is connected with a code, it must be stated if the code is valid (*disco:isValid*) and the code must be stated (*skos:notation*)
  - severity level: ERROR

## 2.70 Recommended Properties<sup>S</sup>

Which properties are not necessarily required but recommended within a particular context<sup>77</sup>.

- ***DATA-CUBE-C-RECOMMENDED-PROPERTIES-01***:
  - severity level: INFO
- ***DCAT-C-RECOMMENDED-PROPERTIES-01***:
  - severity level: INFO
- ***DISCO-C-RECOMMENDED-PROPERTIES-01***: The property *skos:notation* is not mandatory for *disco:Variables*, but recommended to indicate variable names.
  - severity level: INFO
- ***PHDD-C-RECOMMENDED-PROPERTIES-01***:
  - severity level: INFO
- ***SKOS-C-RECOMMENDED-PROPERTIES-01***:
  - severity level: INFO
- ***XKOS-C-RECOMMENDED-PROPERTIES-01***:
  - severity level: INFO

---

<sup>77</sup> *R-72-RECOMMENDED-PROPERTIES*



### 2.71 Handle RDF Collections<sup>C</sup>

Examples of the *Handle RDF Collections*<sup>78</sup> constraint are: a collection must have a specific size; the first/last element of a given list must be a specific literal; the elements of collections are compared; are collections identical?; actions on RDF lists<sup>79</sup>; the 2. list element must be equal to 'XXX'; does the list have more than 10 elements?

- **DISCO-C-HANDLE-RDF-COLLECTIONS-01**: Have comparable *disco:Variables* the same number of codes in their code lists?
  - severity level: INFO
- **DISCO-C-HANDLE-RDF-COLLECTIONS-02**: Does the actual number of *disco:Variables* within an (un)ordered collection of a given *disco:LogicalDataSet* match the expected number?
  - severity level: INFO

### 2.72 Value is Valid for Datatype<sup>B</sup>

Make sure that a value is valid for its datatype. It has to be ensured, e.g., that a date is really a date, or that a *xsd:nonNegativeInteger* value is not negative.

- **DISCO-C-VALUE-IS-VALID-FOR-DATATYPE-01**: Check if all literal values of properties used within the *Disco* context of the datatype *xsd:date* (e.g. *disco:startDate*, *disco:endDate*, *dcterms:date*) are really of the datatype *xsd:date*.
  - severity level: ERROR
- **DISCO-C-VALUE-IS-VALID-FOR-DATATYPE-02**: Frequencies (*disco:frequency*) cannot be negative, i.e., must correspond to the XML Schema datatype *xsd:nonNegativeInteger*.
  - severity level: ERROR
- **DATA-CUBE-C-VALUE-IS-VALID-FOR-DATATYPE-01**: Datatype consistency (*IC-0* [3]) - The RDF graph must be consistent under RDF D-entailment using a datatype map containing all the datatypes used within the graph.
  - severity level: ERROR

### 2.73 Use Sub-Super Relations in Validation<sup>B</sup>

80

The validation of instances data (direct or indirect) exploits the sub-class or sub-property link in a given ontology. This validation can indicate when the data is verbose (redundant) or expressed at a too general level, and could be improved. If *dcterms:date* and one of its sub-properties *dcterms:created* or *dcterms:issued* are present, e.g., check that the value in *dcterms:date* is not redundant with *dcterms:created* or *dcterms:issued* for ingestion.

<sup>78</sup> *R-120-HANDLE-RDF-COLLECTIONS*

<sup>79</sup> See <http://www.snee.com/bobdc.blog/2014/04/rdf-lists-and-sparql.html>

<sup>80</sup> *R-224-USE-SUB-SUPER-RELATIONS-IN-VALIDATION*

- ***DISCO-C-USE-SUB-SUPER-RELATIONS-IN-VALIDATION-01***:  
If one or more *dterms:coverage* properties are present, suggest the use of one of its sub-properties *dterms:spatial* or *dterms:temporal*.
  - severity level: INFO
- ***DISCO-C-USE-SUB-SUPER-RELATIONS-IN-VALIDATION-02***:  
If the *dterms:contributor* property is present, suggest the use of one of its sub-properties, e.g. *disco:fundedBy*.
  - severity level: INFO

## 2.74 Cardinality Shortcuts\*<sup>B</sup>

In most library applications, cardinality shortcuts tend to appear in pairs, with repeatable/non-repeatable establishing maximum cardinality and optional/mandatory establishing minimum cardinality. These are shortcuts for more detailed *cardinality restrictions*.

## 2.75 Aggregation<sup>C</sup>

Some constraints require aggregating multiple values, especially via *COUNT*, *MIN* and *MAX*.

- ***DISCO-C-AGGREGATION-01***: calculate the number of theoretical concepts in the thematic classification of a given study.
  - severity level: INFO
- ***DISCO-C-AGGREGATION-02***: calculate the number of variables of a data set.
  - severity level: INFO
- ***DISCO-C-AGGREGATION-03***: calculate the number of questions in a given questionnaire.
  - severity level: INFO
- ***DISCO-C-AGGREGATION-04***: the number of codes of a given variable must be below a maximum value.
  - severity level: INFO
- ***DISCO-C-AGGREGATION-05***: the number of questions of a given questionnaire must exactly be a given value.
  - severity level: INFO
- ***DISCO-C-AGGREGATION-06***: the sum of percentages of all codes of a given variable must be 100.
  - severity level: INFO
- ***DISCO-C-AGGREGATION-07***: the absolute frequency of all valid codes of a given variable must be equal to a given value.
  - severity level: INFO

## 2.76 Provenance<sup>S</sup>

- **DISCO-C-PROVENANCE-01**: Series should have provenance information (*dcterms:provenance*).
  - severity level: INFO
- **DISCO-C-PROVENANCE-02**: Studies should have provenance information (*dcterms:provenance*).
  - severity level: INFO
- **DISCO-C-PROVENANCE-03**: Data sets should have provenance information (*dcterms:provenance*).
  - severity level: INFO
- **DISCO-C-PROVENANCE-04**: Data files should have provenance information (*dcterms:provenance*).
  - severity level: INFO

## 3 Data Model Specific RDF Constraint Types

### 3.1 Comparison<sup>C</sup>

- **DISCO-C-COMPARISON-VARIABLES-01**: are compared variables represented in a compatible way, i.e. are the variables' code lists theoretically comparable?
  - severity level: WARNING
- **DISCO-C-COMPARISON-VARIABLES-02**: are variable definitions (*dcterms:description*) available for each variable (*disco:Variable*) to compare?
  - severity level: ERROR
- **DISCO-C-COMPARISON-VARIABLES-03**: are code lists structured properly for each variable (*disco:Variable*) to compare?
  - severity level: ERROR
- **DISCO-C-COMPARISON-VARIABLES-04**: is for each code (for each variable (*disco:Variable*) to compare) an associated category (a human-readable label) specified?
  - severity level: INFO
- **DISCO-C-COMPARISON-VARIABLES-05**: each (*disco:Variable*) to compare must be present.
  - severity level: ERROR

### 3.2 Data Model Consistency<sup>C</sup>

Is the data consistent with the intended semantics of the data model? Such validation rules ensure the integrity of the data according to the data model.

- **DISCO-C-DATA-MODEL-CONSISTENCY-01**: Codes (*skos:Concept*) are ordered and therefore have fixed positions in an ordered collection (*skos:OrderedCollection*) within a variable representation. The cumulative percentage of the current code is the cumulative percentage of the previous code (*disco:cumulativePercentage*) plus the percentage value (*disco:percentage*) of the current code.

- severity level: ERROR
- **DISCO-C-DATA-MODEL-CONSISTENCY-02:** The cumulative percentage (*disco:cumulativePercentage*) of the last code must be 100.
  - severity level: ERROR
- **DISCO-C-DATA-MODEL-CONSISTENCY-03:** The number of valid cases (*disco:SummaryStatistics* of the type (*disco:summaryStatisticType*) *ddicv-sumstats:ValidCases*) for a particular variable must exactly be the sum of all frequencies of all valid cases (*disco:inValid* of *skos:Concept* is true).
  - severity level: ERROR
- **DISCO-C-DATA-MODEL-CONSISTENCY-04:** The number of invalid cases (*disco:SummaryStatistics* of the type (*disco:summaryStatisticType*) *ddicv-sumstats:InvalidCases*) for a particular variable must exactly be the sum of all frequencies of all invalid cases (*disco:inValid* of *skos:Concept* is false).
  - severity level: ERROR
- **DISCO-C-DATA-MODEL-CONSISTENCY-05:** The total number of cases (*rdf:value* of the *disco:SummaryStatistics* resource of the type (*disco:summaryStatisticType*) *ddicv-sumstats:NumberOfCases*) for a particular variable must exactly be the number of valid cases plus the number of invalid cases.
  - severity level: ERROR
- **DISCO-C-DATA-MODEL-CONSISTENCY-06:** Some summary statistics types can only be calculated for given variable types. It is not possible to compute minimum values for string variables.
  - severity level: ERROR
- **DISCO-C-DATA-MODEL-CONSISTENCY-07:** Some summary statistics types can only be calculated for given variable types. It is not possible to compute mean values for categorical variables, only for metric variables.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-01:** Only attributes may be optional (*IC-6* [3]) - The only components of a *qb:DataStructureDefinition* that may be marked as optional, using *qb:componentRequired* are attributes.
  - severity level: WARNING
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-02:** Slice Keys consistent with DSD (*IC-8* [3]) - Every *qb:componentProperty* on a *qb:SliceKey* must also be declared as a *qb:component* of the associated *qb:DataStructureDefinition*.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-03:** Slice dimensions complete (*IC-10* [3]) - Every *qb:Slice* must have a value for every dimension declared in its *qb:sliceStructure*.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-04:** All dimensions required (*IC-11* [3]) - Every *qb:Observation* has a value for each dimension declared in its associated *qb:DataStructureDefinition*.
  - severity level: ERROR

- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-05:** No duplicate observations (*IC-12* [3]) - No two *qb:Observations* in the same *qb:DataSet* may have the same value for all dimensions.
  - severity level: WARNING
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-06:** Required attributes (*IC-13* [3]) - Every *qb:Observation* has a value for each declared attribute that is marked as required.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-07:** All measures present (*IC-14* [3]) - In a *qb:DataSet* which does not use a Measure dimension then each individual *qb:Observation* must have a value for every declared measure.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-08:** Measure dimension consistent (*IC-15* [3]) - In a *qb:DataSet* which uses a Measure dimension then each *qb:Observation* must have a value for the measure corresponding to its given *qb:measureType*.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-09:** Single measure on measure dimension observation (*IC-16* [3]) - In a *qb:DataSet* which uses a Measure dimension then each *qb:Observation* must only have a value for one measure (by *IC-15* this will be the measure corresponding to its *qb:measureType*).
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-10:** All measures present in measures dimension cube (*IC-17* [3]) - In a *qb:DataSet* which uses a Measure dimension then if there is a *Observation* for some combination of non-measure dimensions then there must be other *Observations* with the same non-measure dimension values for each of the declared measures.
  - severity level: ERROR
- **DATA-CUBE-C-DATA-MODEL-CONSISTENCY-11:** Consistent data set links (*IC-18* [3]) - 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*.
  - severity level: WARNING
- **SKOS-C-DATA-MODEL-CONSISTENCY-01<sup>81</sup>:** Relation Clashes: Covers condition S27 from the SKOS reference document, that has not been defined formally.
  - Implementation: In a first step, all pairs of concepts are found that are associatively connected, using a SPARQL query. In the second step, a graph is created, containing only hierarchically related concepts and the respective relations. For each concept pair from the first step, we check for a path in the graph from step two. If such a path is found, a clash has been identified and the causing concepts are returned.

<sup>81</sup> Corresponds to qSKOS Quality Issues - SKOS Semi-Formal Consistency Issues - Relation Clashes

- Severity level: INFO
- **SKOS-C-DATA-MODEL-CONSISTENCY-02<sup>82</sup>**: Mapping Clashes: Covers condition S46 from the SKOS reference document, that has not been defined formally.
  - Implementation: Can be solved by issuing a SPARQL query.
  - Severity level: INFO
- **SKOS-C-DATA-MODEL-CONSISTENCY-03<sup>83</sup>**: Mapping Relations Misuse: According to the SKOS reference documentation, mapping relations (e.g., *skos:broadMatch* or *skos:relatedMatch*) should be asserted to concepts being members of different concept schemes. This check finds concepts that are related by a mapping property and are either members of the same concept scheme or members of no concept scheme at all.
  - Severity level: INFO

### 3.3 Structure<sup>C</sup>

SKOS is based on RDF, which is a graph-based data model. Therefore we can concentrate on the vocabulary's graph-based structure for assessing the quality of SKOS vocabularies and apply graph- and network-analysis techniques.

- **DISCO-C-STRUCTURE-01**: there must be exactly one root in the hierarchy of DDI concepts.
  - severity level: ERROR
- **DATA-CUBE-C-STRUCTURE-01**: Codes from hierarchy (*IC-20* [3])
  - If a dimension property has a *qb:HierarchicalCodeList* with a non-blank *qb:parentChildProperty* then the value of that dimension property on every *qb:Observation* must be reachable from a root of the hierarchy using zero or more hops along the *qb:parentChildProperty* links.
    - severity level: ERROR
- **DATA-CUBE-C-STRUCTURE-02**: Codes from hierarchy (inverse) (*IC-21* [3])
  - If a dimension property has a *qb:HierarchicalCodeList* with an inverse *qb:parentChildProperty* then the value of that dimension property on every *qb:Observation* must be reachable from a root of the hierarchy using zero or more hops along the inverse *qb:parentChildProperty* links.
    - severity level: ERROR
- **SKOS-C-STRUCTURE-01<sup>84</sup>**: Orphan Concepts: An orphan concept is a concept without any associative or hierarchical relations. It might have attached literals like e.g., labels, but is not connected to any other resource, lacking valuable context information. A controlled vocabulary that contains many orphan concepts is less usable for search and retrieval use cases, because, e.g., no hierarchical query expansion can be performed on search terms to find documents with more general content.

<sup>82</sup> Corresponds to qSKOS Quality Issues - SKOS Semi-Formal Consistency Issues - Mapping Clashes

<sup>83</sup> Corresponds to qSKOS Quality Issues - SKOS Semi-Formal Consistency Issues - Mapping Relations Misuse

<sup>84</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Orphan Concepts

- Implementation: Iteration over all concepts in the vocabulary and returning that don't have associated resources using (sub-properties of) *skos:semanticRelation*.
- Severity level: WARNING
- **SKOS-C-STRUCTURE-02<sup>85</sup>**: Disconnected Concept Clusters: Checking the connectivity of the graph, it is possible to identify all weakly connected components. These datasets form "islands" in the vocabulary and might be caused by incomplete data acquisition, "forgotten" test data, outdated terms and the like.
  - Implementation: Creation of an undirected graph that includes all non-orphan concepts as nodes and all semantic relations as edges. Tarjan's algorithm then finds and returns all weakly connected components.
  - Severity level: INFO
- **SKOS-C-STRUCTURE-03<sup>86</sup>**: Cyclic Hierarchical Relations: Although perfectly consistent with the SKOS data model, cyclic relations may reveal a logical problem in the thesaurus. Consider the following example: "decision" → "problem resolution" → "problem" (→ "decision": here the cycle is closed). The concepts are connected using *skos:broader* relationships (indicated with "→"). Due to the fact that a thesaurus is in many cases a product of consensus between the contributors (or just the decision of one dedicated thesaurus manager), it will be almost impossible to automatically resolve the cycle (i.e. deleting an edge).
  - Implementation: Construction of a graph having all concepts as nodes and the set of edges being *skos:broader* relations.
  - Severity level: WARNING
- **SKOS-C-STRUCTURE-04<sup>87</sup>**: Valueless Associative Relations: Two concepts are sibling, but also connected by an associative relation. In that context, the associative relation is not necessary. See ISO\_DIS\_25964-1, 11.3.2.2
  - Implementation: Identification of all pairs of concepts that have the same broader or narrower concepts, i.e. they are "sibling terms". All siblings that are related by a *skos:related* property are returned.
  - Severity level: INFO
- **SKOS-C-STRUCTURE-05<sup>88</sup>**: Solely Transitively Related Concepts: *skos:broaderTransitive* and *skos:narrowerTransitive* are, according to the SKOS reference document, "not used to make assertions", so they should not be the only relations hierarchically relating two concepts.
  - Implementation: Identification of all concept pairs that are related by *skos:broaderTransitive* or *skos:narrowerTransitive* properties but not by their *skos:broader* and *skos:narrower* subproperties.

<sup>85</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Disconnected Concept Clusters

<sup>86</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Cyclic Hierarchical Relations

<sup>87</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Valueless Associative Relations

<sup>88</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Solely Transitively Related Concepts

- Severity level: INFO
- **SKOS-C-STRUCTURE-06<sup>89</sup>**: Unidirectionally Related Concepts: Reciprocal relations (e.g., *broader/narrower*, *related*, *hasTopConcept/topConceptOf*) should be included in the controlled vocabularies to achieve better search results using SPARQL in systems without reasoner support.
  - Implementation: This issue is checked without inference of *owl:inverseOf* properties. We iterate over all triples and check for each property if an inverse property is defined in the SKOS ontology and if the respective statement using this property is included in the vocabulary. If not, the resources associated with this property are returned.
  - Severity level: INFO
- **SKOS-C-STRUCTURE-07<sup>90</sup>**: Omitted Top Concepts: A vocabulary should provide "entry points" to the data to provide "efficient access" (SKOS primer) and guidance for human users.
  - Implementation: For every ConceptScheme in the controlled vocabulary, a SPARQL query is issued finding resources that are associated with this ConceptScheme by one of the properties *skos:hasTopConcept* or *skos:topConceptOf*. Top concepts are also concepts having no broader concept.
  - Severity level: WARNING
- **SKOS-C-STRUCTURE-08<sup>91</sup>**: Top Concepts Having Broader Concepts: Concepts "internal to the tree" should not be indicated as top concepts.
  - Implementation: A SPARQL query finds all top concepts (being defined by one of the properties *skos:hasTopConcept* or *skos:topConceptOf*) having associated a broader concept.
  - Severity level: ERROR
- **SKOS-C-STRUCTURE-09<sup>92</sup>**: Hierarchical Redundancy: As stated in the SKOS reference document, *skos:broader* and *skos:narrower* are not transitive properties. However, they are sub-properties of *skos:broaderTransitive* and *skos:narrowerTransitive* which enables inference of a "transitive closure". This, in fact, leaves it up to the user to interpret whether a vocabulary's hierarchical structure is seen as transitive or not. In the former case, this check can be useful. It finds pairs of concepts (A,B) that are directly hierarchically related but there also exists an hierarchical path through a concept C that connects A and B.
  - Severity level: INFO
- **SKOS-C-STRUCTURE-10<sup>93</sup>**: Reflexive Relations: Concepts related to themselves.
  - Severity level: WARNING

<sup>89</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Unidirectionally Related Concepts

<sup>90</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Omitted Top Concepts

<sup>91</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Top Concepts Having Broader Concepts

<sup>92</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Hierarchical Redundancy

<sup>93</sup> Corresponds to qSKOS Quality Issues - Structural Issues - Reflexive Relations



### 3.4 Labeling and Documentation<sup>S</sup>

- **DISCO-C-LABELING-AND-DOCUMENTATION-01**: Series should be described (*dcterms:description*).
  - severity level: INFO
- **DISCO-C-LABELING-AND-DOCUMENTATION-02**: Studies should be described (*dcterms:description*).
  - severity level: INFO
- **DISCO-C-LABELING-AND-DOCUMENTATION-03**: Data sets should be described (*dcterms:description*).
  - severity level: INFO
- **DISCO-C-LABELING-AND-DOCUMENTATION-04**: Data files should be described (*dcterms:description*).
  - severity level: INFO
- **DISCO-C-LABELING-AND-DOCUMENTATION-05**: Instruments should be described (*dcterms:description*).
  - severity level: INFO
- **DISCO-C-LABELING-AND-DOCUMENTATION-06**: Variables should be described (*dcterms:description*).
  - severity level: INFO
- **SKOS-C-LABELING-AND-DOCUMENTATION-01<sup>94</sup>**: Undocumented Concepts: The SKOS standard defines a number of properties useful for documenting the meaning of the concepts in a thesaurus also in a human-readable form. Intense use of these properties leads to a well-documented thesaurus which should also improve its quality.
  - Implementation: Iteration over all concepts in the vocabulary and find those not using one of *skos:note*, *skos:changeNote*, *skos:definition*, *skos:editorialNote*, *skos:example*, *skos:historyNote*, or *skos:scopeNote*.
  - Severity level: INFO
- **SKOS-C-LABELING-AND-DOCUMENTATION-02<sup>95</sup>**: Overlapping Labels: This is a generalization of a recommendation in the SKOS primer, that “no two concepts have the same preferred lexical label in a given language when they belong to the same concept scheme”. This could indicate missing disambiguation information and thus lead to problems in autocompletion application.
  - Severity level: INFO
- **SKOS-C-LABELING-AND-DOCUMENTATION-03<sup>96</sup>**: Missing Labels: To make the vocabulary more convenient for humans to use, instances of SKOS classes (Concept, ConceptScheme, Collection) should be labeled using e.g., *skos:prefLabel*, *altLabel*, *rdfs:label*, *dc:title*.

<sup>94</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Undocumented Concepts

<sup>95</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Overlapping Labels

<sup>96</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Missing Labels

- Severity level: INFO
- ***SKOS-C-LABELING-AND-DOCUMENTATION-04***<sup>97</sup>: Unprintable Characters in Labels: *pref/alt/hiddenlabels* contain characters that are not alphanumeric characters or blanks.
  - Severity level: INFO
- ***SKOS-C-LABELING-AND-DOCUMENTATION-05***<sup>98</sup>: Empty Labels: Labels also need to contain textual information to be useful, thus we find all SKOS labels with length 0 (after removing whitespaces).
  - Severity level: INFO
- ***SKOS-C-LABELING-AND-DOCUMENTATION-06***<sup>99</sup>: Ambiguous Notation References: Concepts within the same concept scheme should not have identical *skos:notation* literals.
  - Severity level: INFO

### 3.5 Vocabulary<sup>B</sup>

Vocabularies should not invent any new terms or use deprecated elements.

- ***DATA-CUBE-C-VOCABULARY-01***
  - Severity level: ERROR
- ***DCAT-C-VOCABULARY-01***
  - Severity level: ERROR
- ***DISCO-C-VOCABULARY-01***
  - Severity level: ERROR
- ***PHDD-C-VOCABULARY-01***
  - Severity level: ERROR
- ***SKOS-C-VOCABULARY-01***<sup>100</sup>: Undefined SKOS Resources: The vocabulary should not invent any new terms within the SKOS namespace or use deprecated SKOS elements.
  - Severity level: ERROR
- ***XKOS-C-VOCABULARY-01***
  - Severity level: ERROR

<sup>97</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Unprintable Characters in Labels

<sup>98</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Empty Labels

<sup>99</sup> Corresponds to qSKOS Quality Issues - Labeling and Documentation Issues - Ambiguous Notation References

<sup>100</sup> Corresponds to qSKOS Quality Issues - Linked Data Specific Issues - Undefined SKOS Resources

### 3.6 HTTP URI Scheme Violation<sup>S</sup>

- ***DISCO-C-HTTP-URI-SCHEME-VIOLATION*<sup>101</sup>**: In the context of Linked Data, we restrict ourselves to using HTTP URIs only and avoid other URI schemes such as URNs and DOIs.
  - Severity level: ERROR
- ***DATA-CUBE-C-HTTP-URI-SCHEME-VIOLATION*<sup>102</sup>**: In the context of Linked Data, we restrict ourselves to using HTTP URIs only and avoid other URI schemes such as URNs and DOIs.
  - Severity level: ERROR
- ***PHDD-C-HTTP-URI-SCHEME-VIOLATION*<sup>103</sup>**: In the context of Linked Data, we restrict ourselves to using HTTP URIs only and avoid other URI schemes such as URNs and DOIs.
  - Severity level: ERROR
- ***SKOS-C-HTTP-URI-SCHEME-VIOLATION*<sup>104</sup>**: In the context of Linked Data, we restrict ourselves to using HTTP URIs only and avoid other URI schemes such as URNs and DOIs.
  - Severity level: ERROR

## 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, Andreas Nolle, Erman Acar, and Kai Eckert. Rdf validation requirements - evaluation and logical underpinning. 2015.
3. Richard Cyganiak and Dave Reynolds. The rdf data cube vocabulary. W3C recommendation, W3C, January 2014.
4. Dimitris Kontokostas, Patrick Westphal, Sören Auer, Sebastian Hellmann, Jens Lehmann, Roland Cornelissen, and Amrapali Zaveri. Test-driven evaluation of linked data quality. In *Proceedings of the 23rd International Conference on World Wide Web, WWW '14*, pages 747–758, Republic and Canton of Geneva, Switzerland, 2014. International World Wide Web Conferences Steering Committee.
5. Markus Krötzsch, Frantisek Simancik, and Ian Horrocks. A description logic primer. *CoRR*, abs/1201.4089, 2012.
6. 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.
7. Michael Schneider. OWL 2 Web Ontology Language RDF-Based Semantics. W3C recommendation, W3C, October 2009.

<sup>101</sup> Corresponds to qSKOS Quality Issues - Linked Data Specific Issues - HTTP URI Scheme Violation

<sup>102</sup> Corresponds to qSKOS Quality Issues - Linked Data Specific Issues - HTTP URI Scheme Violation

<sup>103</sup> Corresponds to qSKOS Quality Issues - Linked Data Specific Issues - HTTP URI Scheme Violation

<sup>104</sup> Corresponds to qSKOS Quality Issues - Linked Data Specific Issues - HTTP URI Scheme Violation