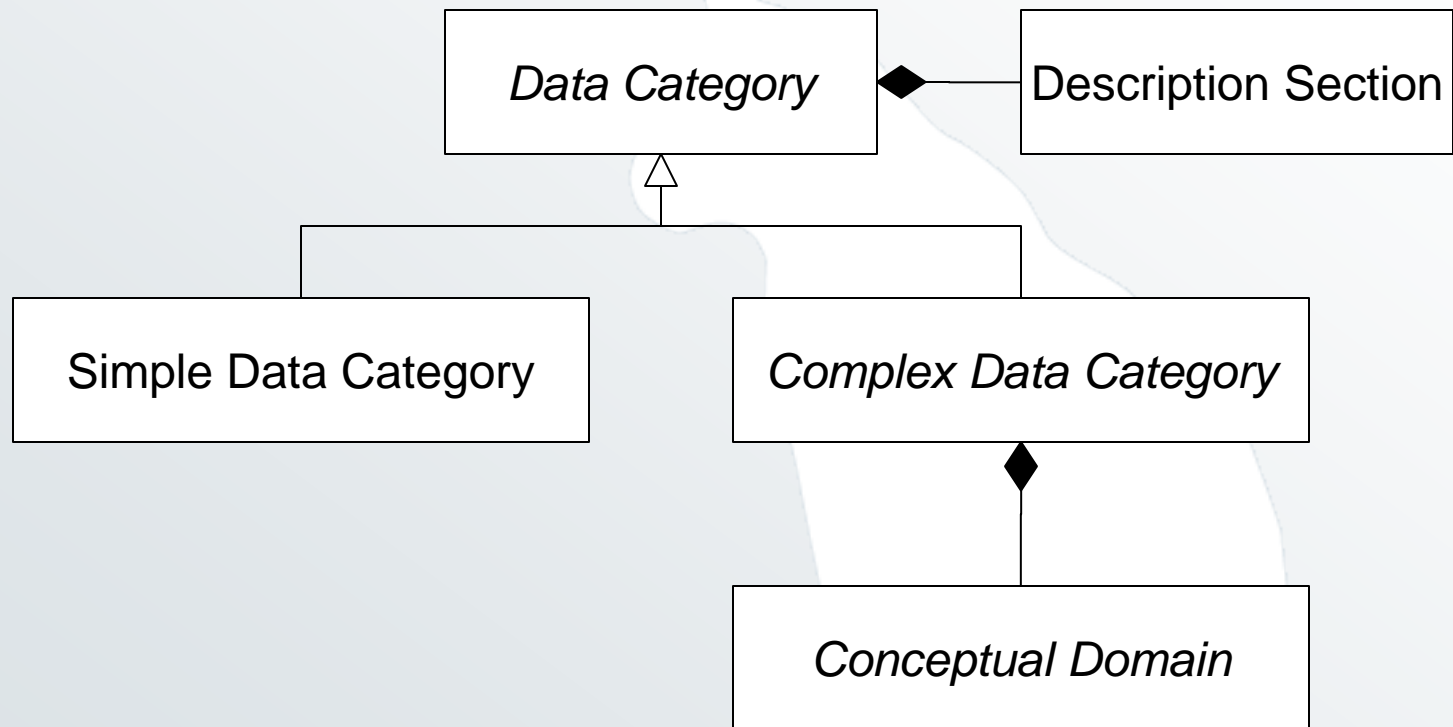


ISOcat tutorial

DCR data model and guidelines

Simple and complex DCs

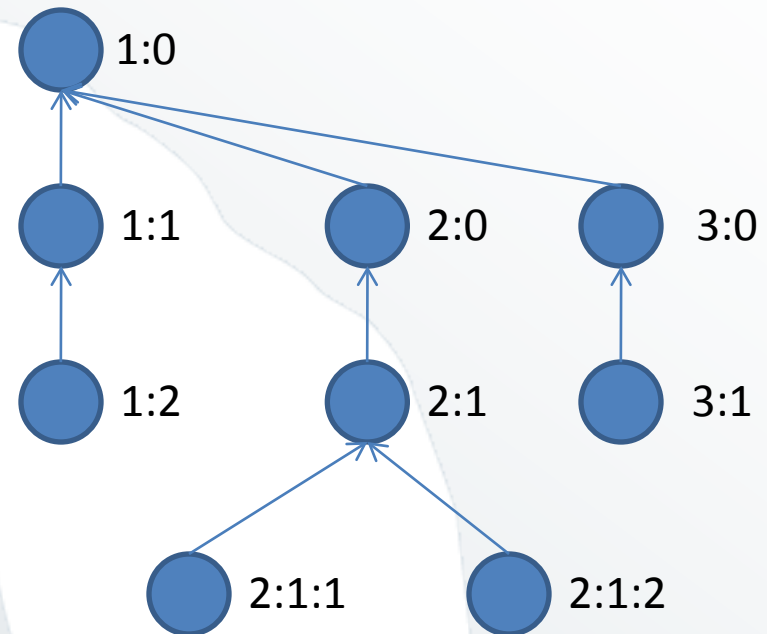


Administrative Information Section

- The mandatory identifier
 - is a mnemonic string used to refer to the data category
 - should be based on a meaningful English word or series of words presented as an alphanumeric character string; for multiword strings, begin with lowercase and express the identifier as one continuous string in camel case with no white space (for instance, */term/*, */normativeAuthorization/*, */preferredTerm/*)
 - maybe used in XML vocabularies and thus must be a valid local part of a qualified name:
 - Cannot start with a number, shouldn't contain any whitespace, ...
- ISOcat warns you when the identifier is invalid and will refuse to save the data category

Administrative Information Section

- The mandatory version
 - used to refine identifier to indicate the version of the data category
 - managed by the system
 - based on branching not on major and minor revisions



Administrative Information Section

- The mandatory data category type
 - complex/open
 - conceptual domain is not restricted to an enumerated set of values
 - complex/constrained
 - conceptual domain is non-enumerated, but is restricted to a constraint specified in a schema-specific language or languages
 - complex/closed
 - conceptual domain is restricted to a set of enumerated simple data categories making up its value domain
 - simple

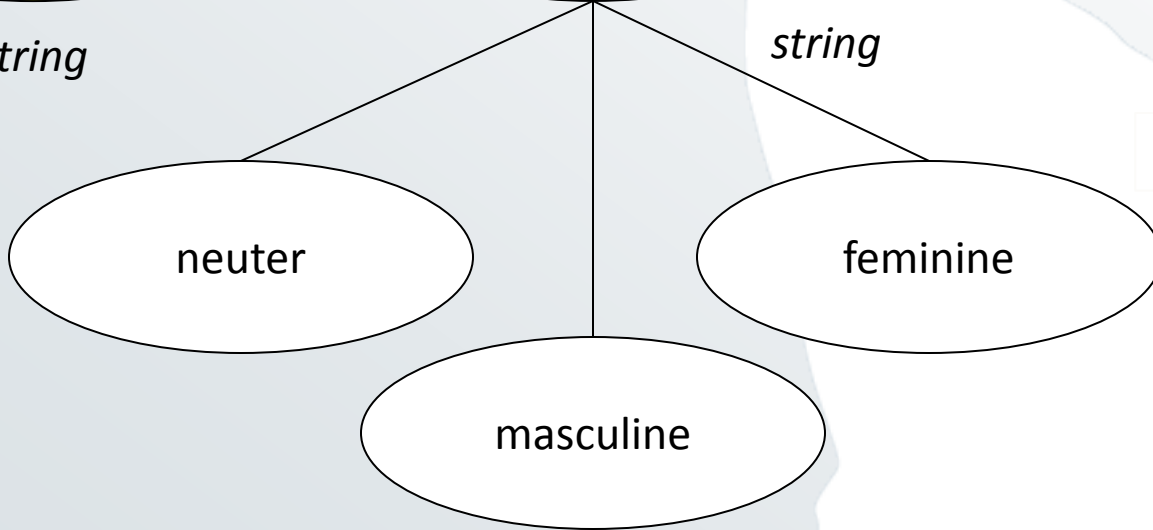
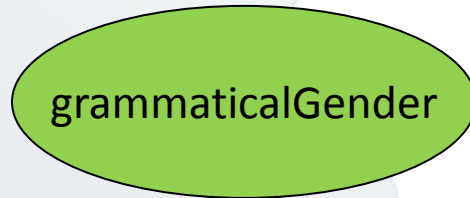
Data category types

complex: open



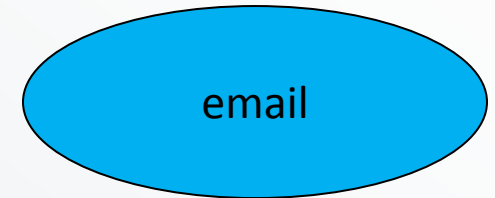
string

closed



string

constrained



string

Constraint: .+@.+

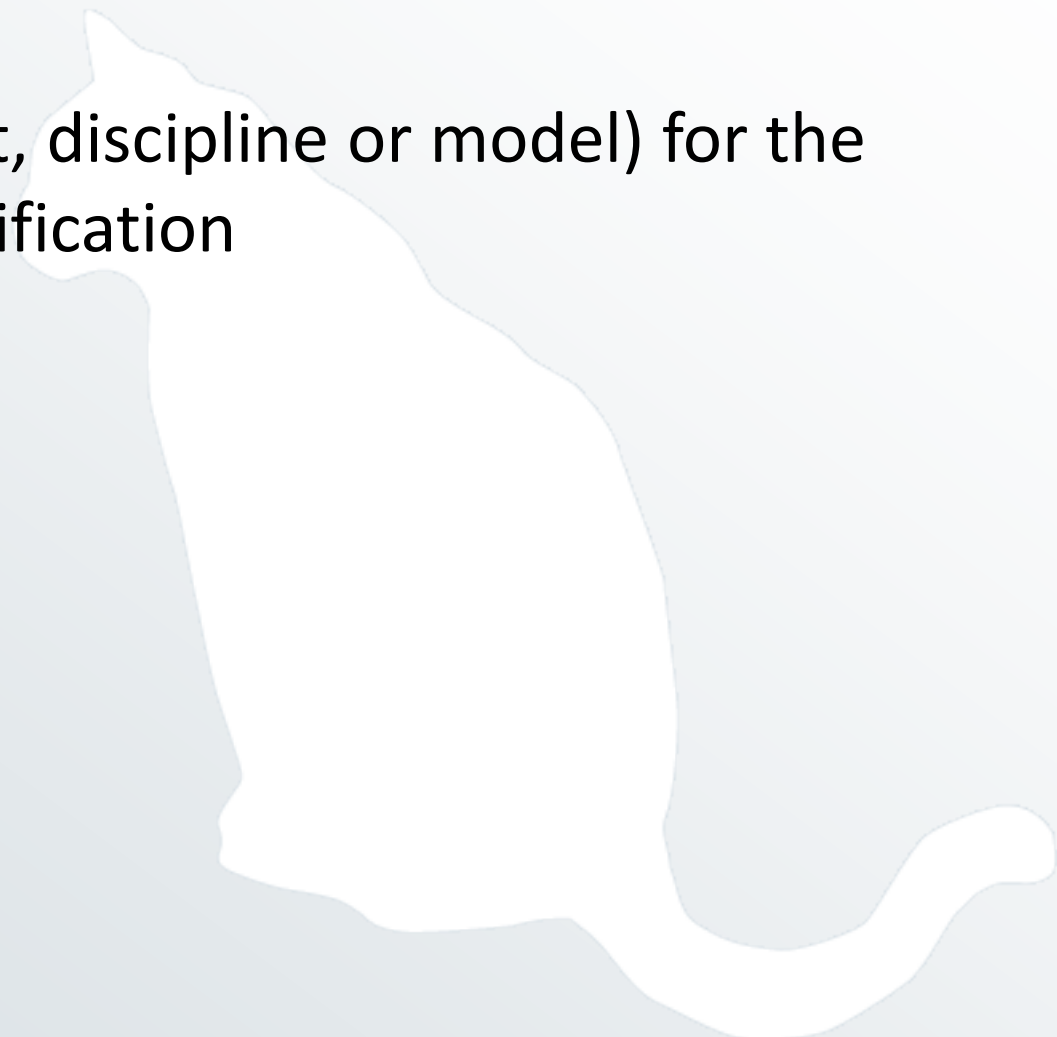
simple:

Administrative Information Section

- The justification
 - a short description justifying why the data category should be included in the registry
 - mandatory for data categories to be standardized; desirable in general
 - even data categories that are common in a given thematic domain may be unfamiliar or ambiguous to users unfamiliar with that domain

Administrative Information Section

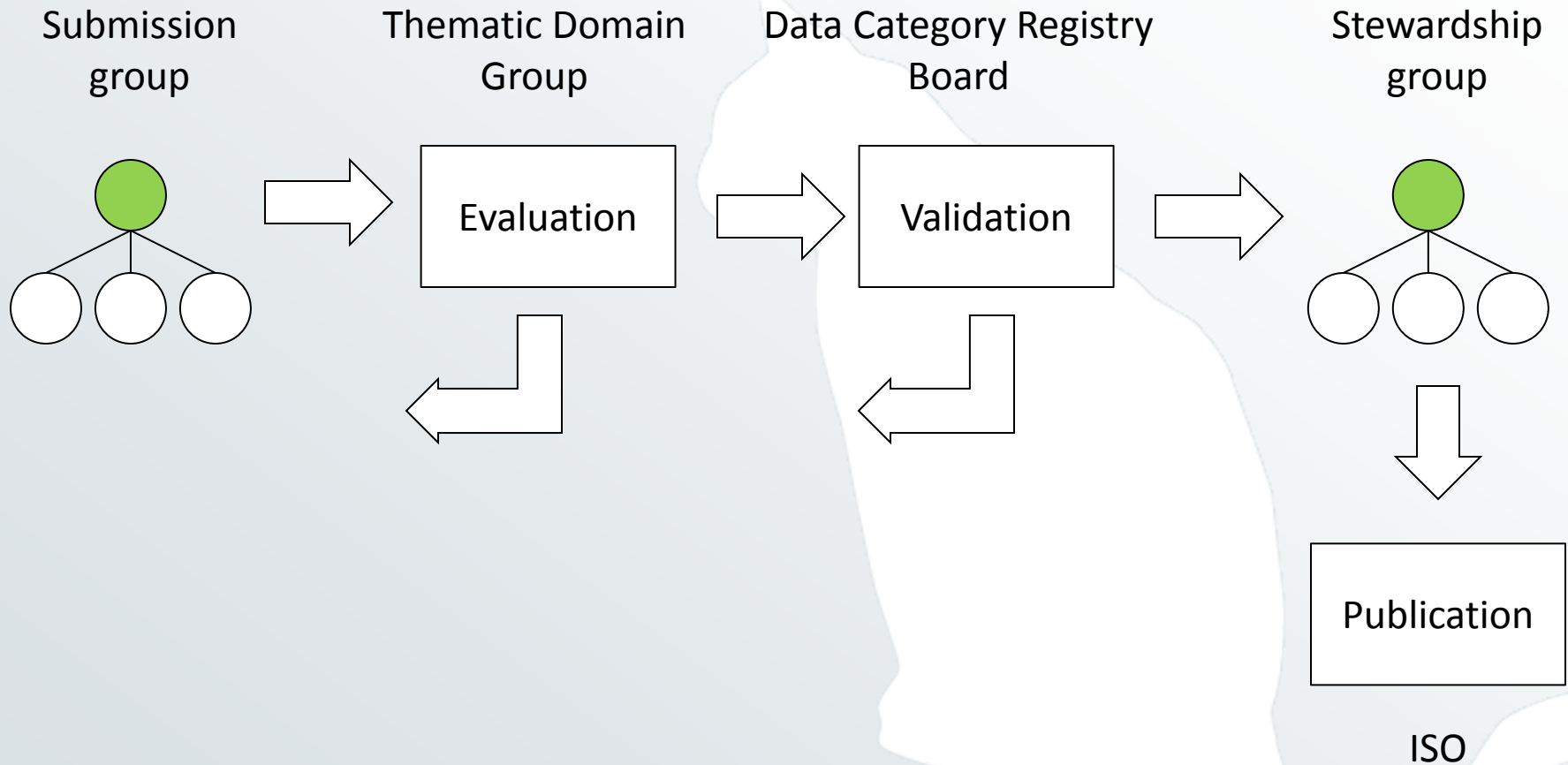
- The origin
 - (document, project, discipline or model) for the data category specification



Administrative Information Section

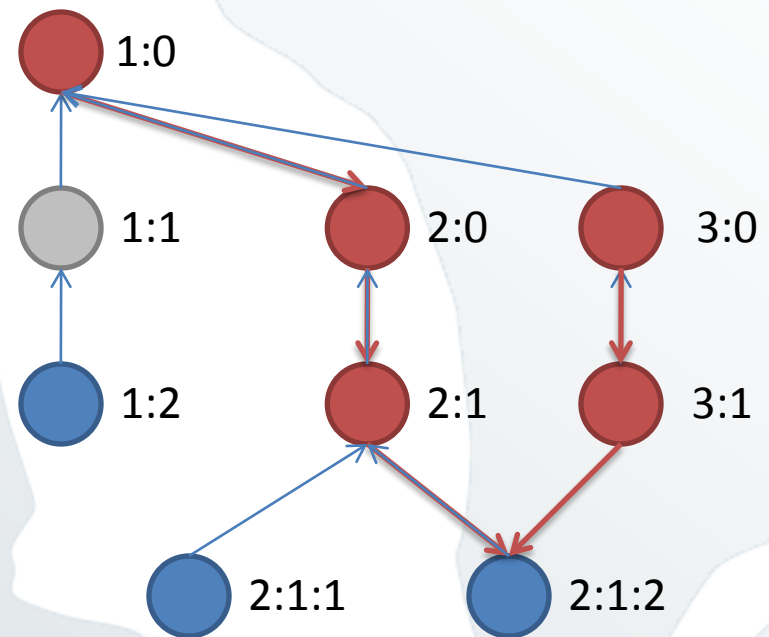
- The mandatory administration status
 - a designation of the status in the administration process for handling registration requests under the stewardship of the DCR Board
 - managed by the system
 - values:
 - private
 - submission
 - pre-evaluation, evaluation, rejected-TDG, accepted-TDG
 - pre-validation, validation, rejected-DCR Board
 - accepted

Standardization process



Administrative Information Section

- The mandatory registration status
 - a designation of the status in the registration life-cycle of an administered item
 - managed by the system
 - values:
 - private
 - candidate
 - standard
 - superseded
 - deprecated



Administrative Information Section

- The effective date
 - the date a data category specification has/will become available to DCR users
- The until date
 - the date a data category specification is no longer effective in the registry; this information is set when the registration status of the data category specification changes to *deprecated* or *superseded*
- ISOcat isn't acting on these dates (yet)

Administrative Information Section

- The explanatory comments
 - descriptive comments about the data category specification
- The unresolved issues
 - a problem that remains unresolved regarding proper documentation of the data category specification

Description Section

- The mandatory profile
 - attribute used to relate the current data category specification to one or several thematic domains treated by ISO/TC 37 (for example, morphosyntax, syntax, metadata, language description, etc.)
 - the value of profile defaults to *Private*
 - submission for standardization requires the selection of at least one thematic domain profile because it is the relevant TDG that is responsible for maintenance of standardized data category specifications
 - if multiple profiles from multiple TDGs are selected one TDG will still be responsible, but the other TDGs will be involved in the harmonization process
 - if a user desires to create a new profile, consult the DCR Board

Profiles and TDGs

- Each Thematic Domain Group (TDG) manages at least one profile
- The following TDGs are active:
 - Metadata
 - Morphosyntax
 - Terminology
- The Language codes profile will be maintained separately and hooked up to ISOcat
- Other TDGs and their profiles are still asleep
- ISOcat will host a forum in the near future, which also includes the functionality to send users, i.e., an owner of a DC or the chair of a TDG, an email

Data Element Name Sections

- used to record names for the data category as used in a given database, format or application
- language independent
- attributes:
 - the mandatory data element name
 - one identifier (word, multi-word unit or (alpha)numeric representation)
 - the mandatory source
 - the database, format or application in which the data element name is used
- ISOcat will future have better support for keeping sources in sync

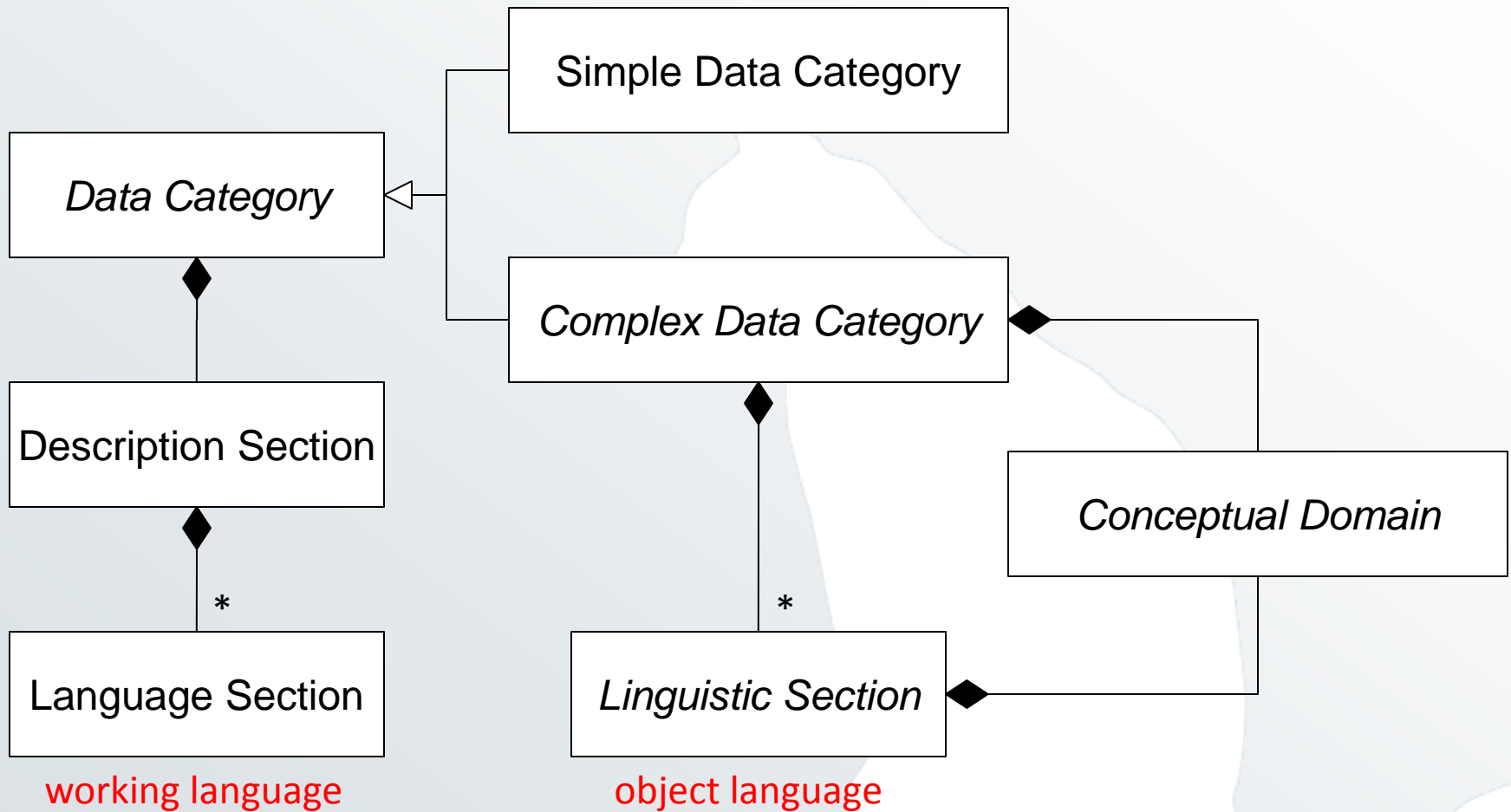
Working and object languages

- Working language:
 - language used to describe objects
- Object language:
 - language being described

You can describe properties of the object language French in the working language Dutch:

In de Franse taal worden vrouwelijke en mannelijk zelfstandige naamwoorden onderscheiden.

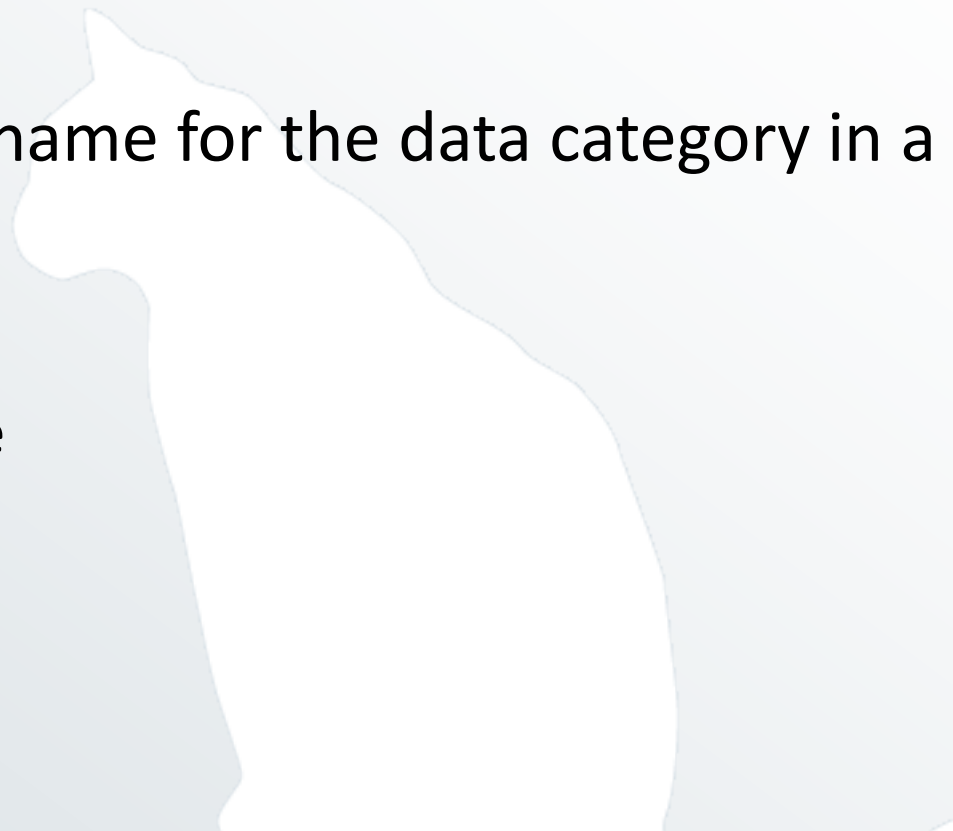
Working and object languages



Language Section

- The English language section is mandatory and has to contain at least one name and one definition
- Additional language sections can be added as needed, and should translate at least the definition of the English language section
- if a user needs an additional language, contact the DCR system administration

Language Section

- The Name Section
 - records a possible name for the data category in a specific language
 - status:
 - standardized name
 - preferred name
 - admitted name
 - deprecated name
 - superseded name
- 

Language Section

- The Definition Section
 - definition of the data category concept associated with the data category, written in the language of the language section
 - attributes:
 - definition:
 - definitive formulation that should be general enough to apply to all thematic domains and implementations of the data category
 - source:
 - from which the definition has been borrowed or adapted
 - note:
 - any additional information about the definition

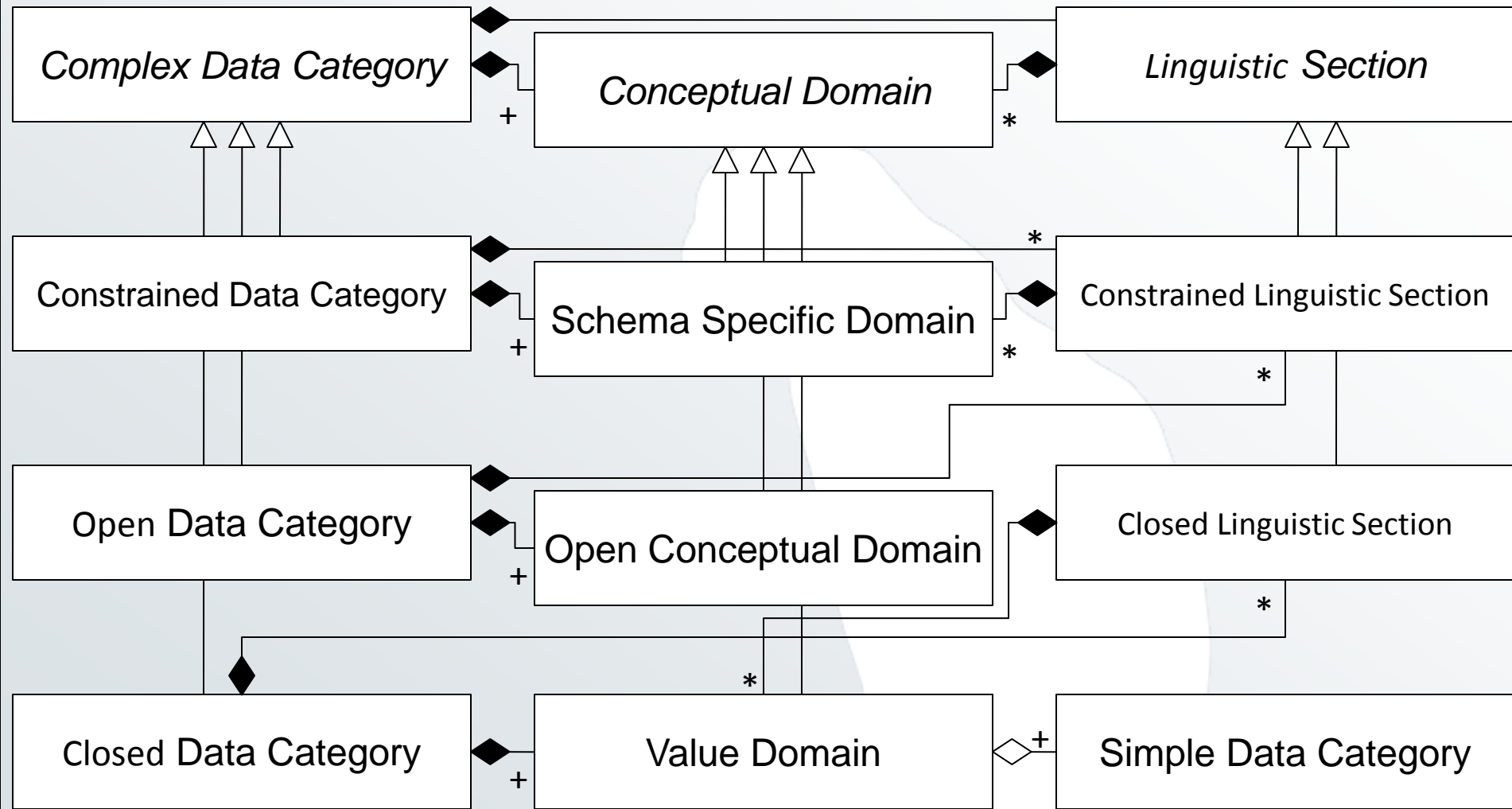
ISO 704 Terminology work — principles and methods

- Intensional definitions:
 - They should consist of a single sentence fragment;
 - They should begin with a meaningful *broader concept*, *either immediately above or at a higher level* of the data category concept being defined;
 - They should list critical and delimiting *characteristic(s) that distinguish the concept from other related concepts*.
- Actual concept systems, such as are implied here by the reference to broader and related concepts, should be modeled in Relation Registries outside the DCR. Furthermore, different domains and communities of practice may differ in their choice of the immediate broader concept, depending upon any given ontological perspective. Harmonized definitions for shared DCs should attempt to choose generic references insofar as possible.

Language Section

- The examples
 - a sample instance reflecting the data category
 - should be limited to those that illustrate the data category in general, excluding language specific usage, which should be documented in a Linguistic Section
 - may be accompanied by the source of the example
- The explanations
 - any additional information about the data category that would not be relevant for a definition (for example, more precise linguistic background concerning the use of the data category)
 - may be accompanied by the source from which the explanation has been borrowed or adapted
- The notes
 - any additional information associated with the data category, excluding technical information that would normally be described in an explanation

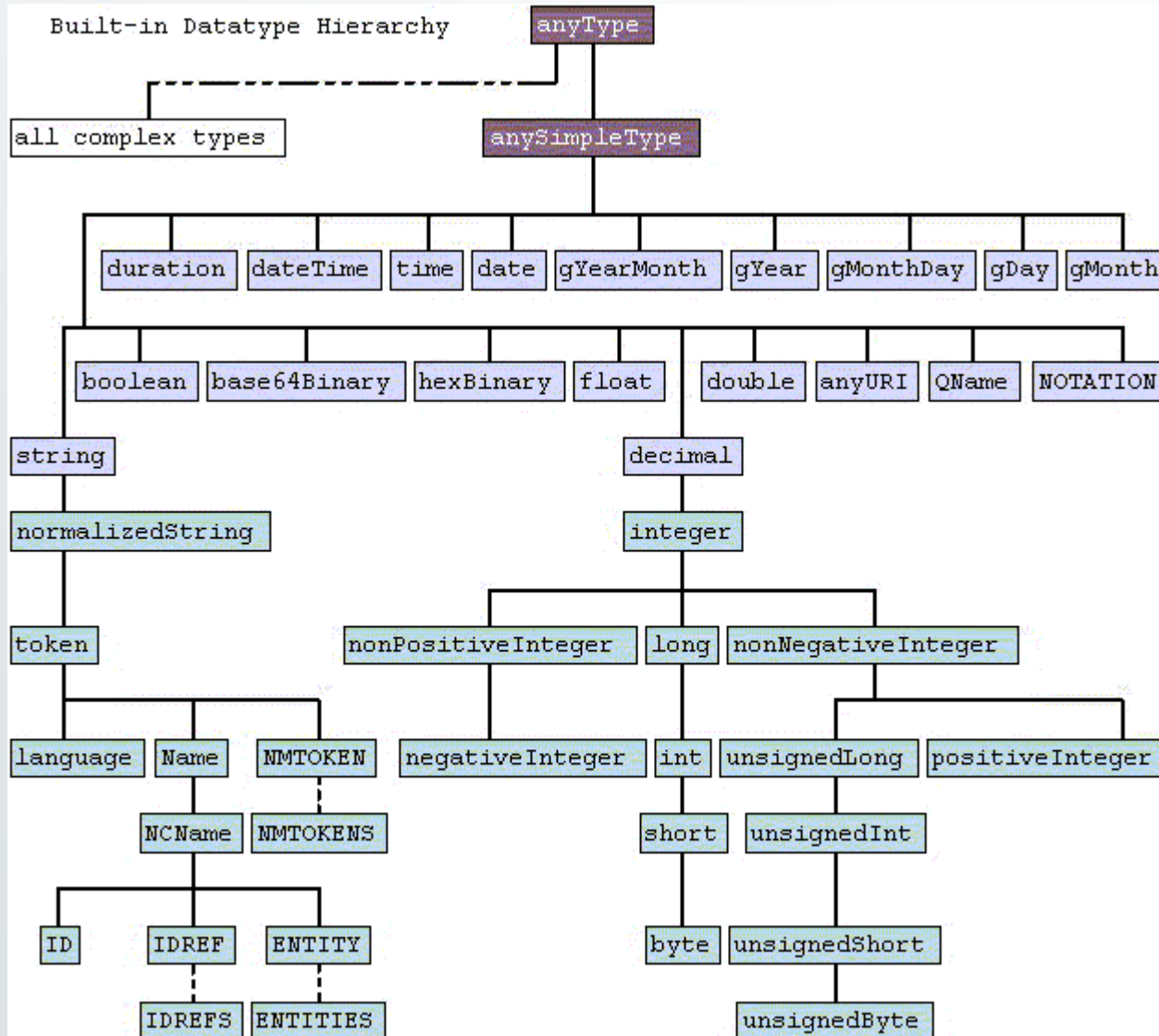
Conceptual domains



Conceptual domains

- The mandatory data type
 - the data type, as defined for W3C XML Schema, of this complex data category
 - the default data type is *string*
- Just the data type is enough for an open data category

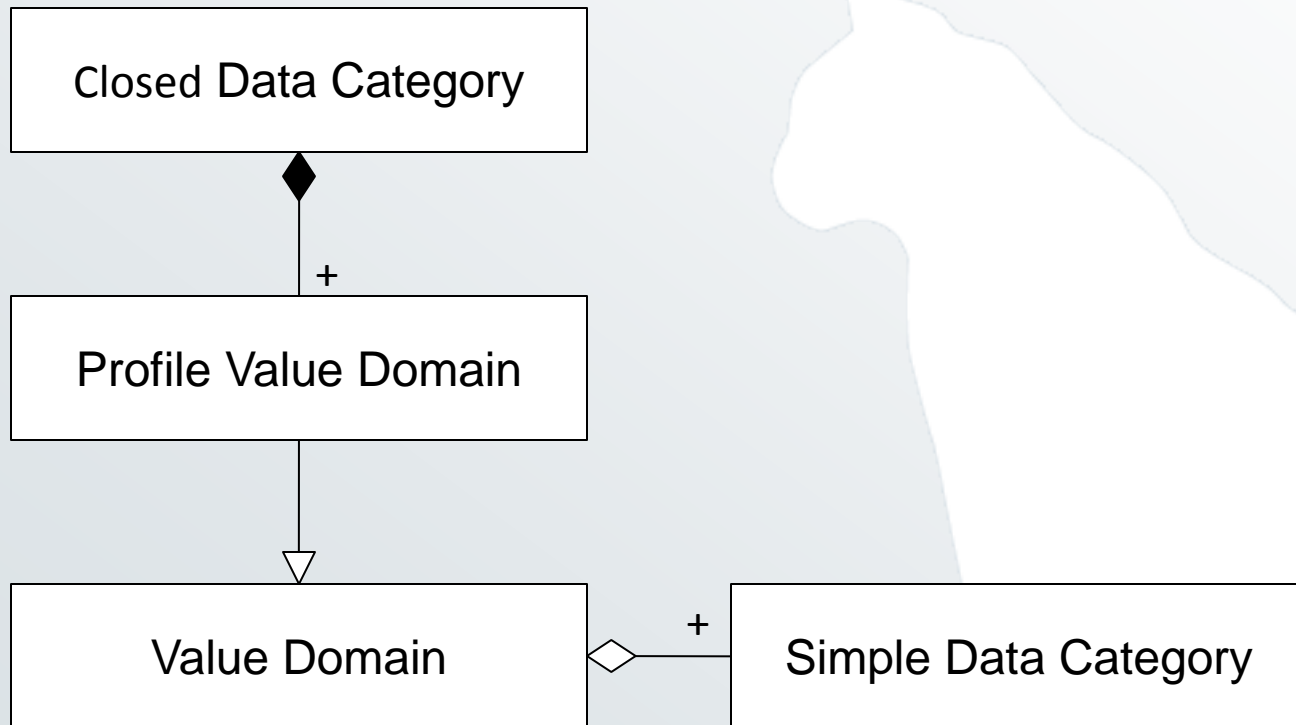
W3X XML Schema data types



Constrained Conceptual Domain

- The mandatory constraint:
 - allows users to express constraints on the possible values of a conceptual domain associated with a given data type in a rule language suitable for the schema in question
 - rule languages currently ‘supported’:
 - Schematron
 - Object Constraint Language
 - Semantic Web Rule Language
 - Relax NG
 - the DCR doesn’t do any interpretation of the rules, this has to be done by the user/TDG/DCR Board
- ISOcat does check if its valid XML (when applicable)
- ISOcat’s DCIF export is somewhat incomplete in this area: constraints expressed in XML are outputted as one text block instead of being an integrated part of the DCIF document

Profile value domains



Profile Value Domain

- set of permissible values for a specific profile
- each value is represented by a simple data category
- the simple data category needs to be a member of the profile

- ISOcat makes an exception for the *Private* profile value domain which can contain any simple data category
- Standardized data categories can't have a *Private* profile value domain
- ISOcat tries to simplify the creation of the value domain by automating the assignment to profile value domains. However, at the moment this is too simplistic and a more advanced editor is needed to fully manage profile value domains.

Example

Data category	Morposyntax	Terminology
<i>/partOfSpeech/</i>	X	X
<i>/adjective/</i>	X	X
<i>/ordinalAdjective/</i>	X	
<i>/participleAdjective/</i>	X	
<i>/qualifierAdjective/</i>	X	
<i>/adposition/</i>	X	X
<i>/circumposition/</i>	X	
<i>/preposition/</i>	X	
<i>/postposition/</i>	X	

Linguistic Sections

- used to specify the behaviour of a complex data category in a specific object language
- Language specific examples, explanations and notes
- Refinement of the conceptual domain:
 - (additional) constraints for open and constrained complex data categories
 - subset value domains for closed complex data categories

Example

- /grammaticalGender/ has as (profile specific) value domain:
 - /*masculine*/
 - /*feminine*/
 - /*neuter*/
- The French linguistic section limits it to:
 - /*masculine*/
 - /*feminine*/

Hierarchical Simple Data Categories

- Simple data categories can be put in a subsumption (is-a) hierarchy
 - allows different levels of granularity in a value domain
 - make large value domains manageable
 - a simple data category can be only a member of one hierarchy, i.e., it can have only one parent

Changes

- Each time you save a data category you're asked to enter a description of what you've changed
- These descriptions are available in the history log of each data category
- The DCIF contains the first (upon time of creation) and the last change description

Checker

- Not all mandatory parts of the specification has to be filled at once
- The check will tell you what is still missing
- Some errors or warnings can only be fixed by issues change requests to a standardized data category
 - Membership of another profile