EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

SFB 833, Project INF:
**Heterogene Forschungsprimärdaten des SFB 833
Repräsentation und Verarbeitung**

NALiDa
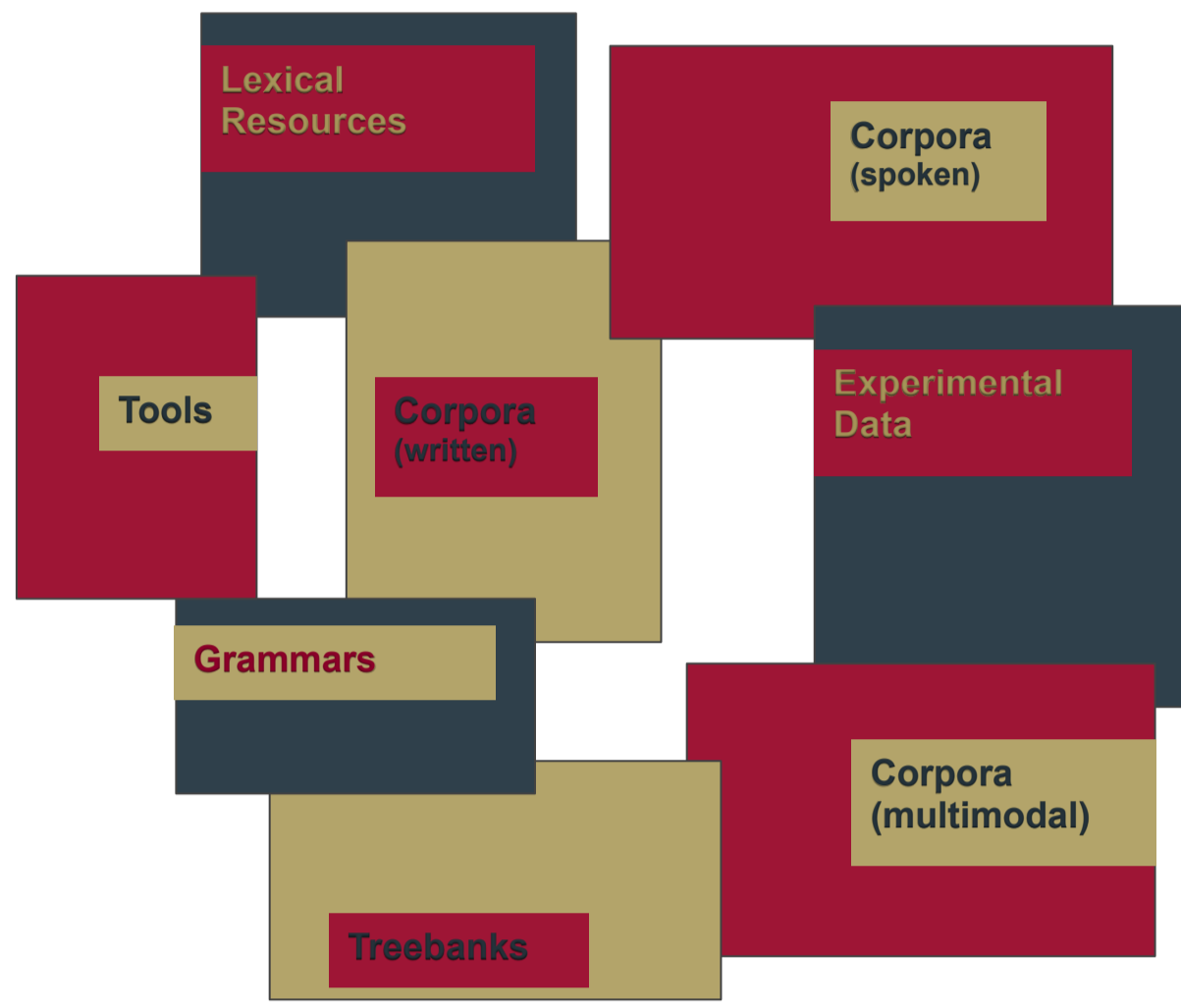Centre for Sustainability of Linguistic Data

# Managing Linguistic Resources by Enriching Their Metadata with Linked Data
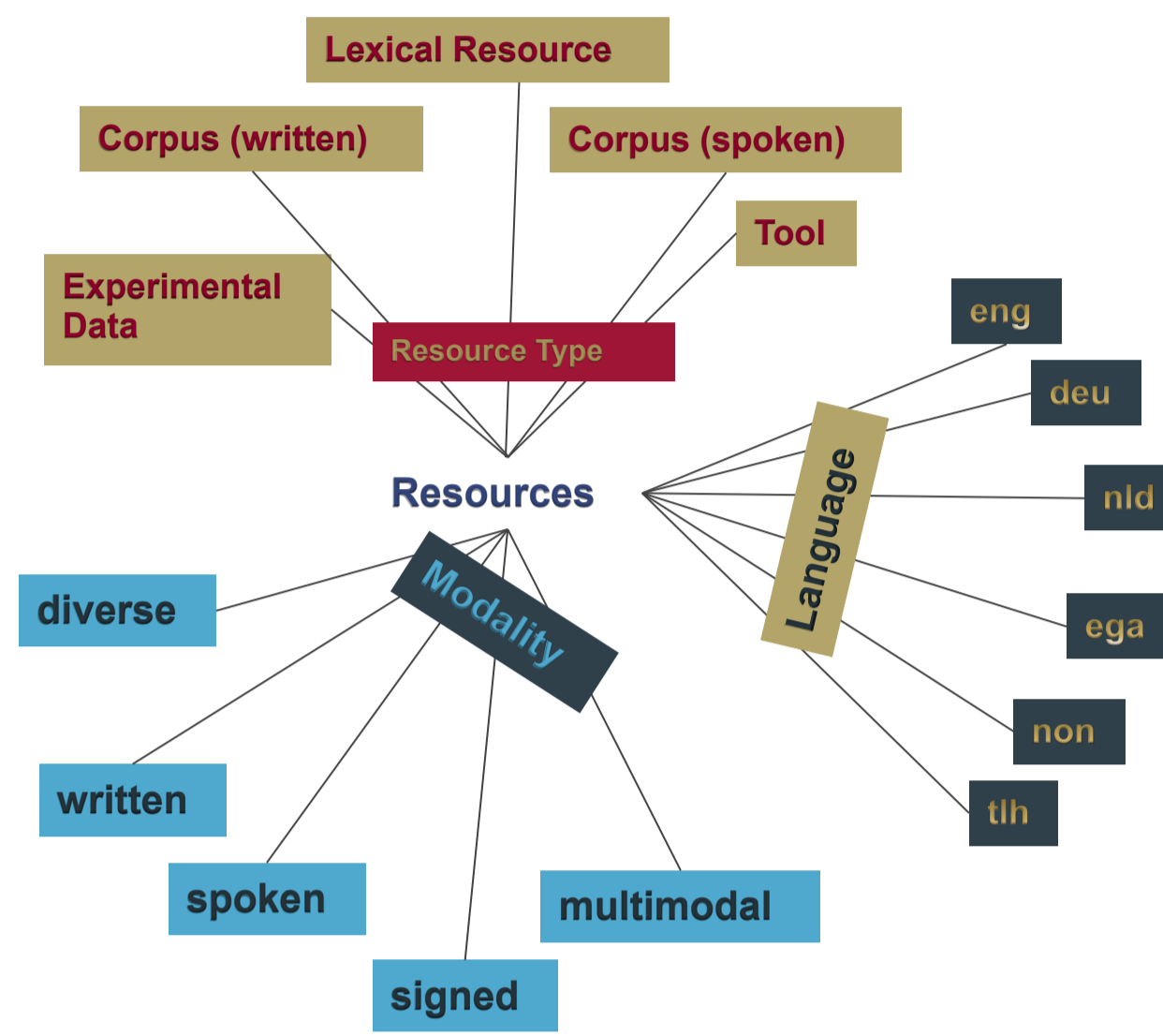
**Christina Hoppermann
Thorsten Trippel
Claus Zinn**

## Sustainable Management of Linguistic Resources

### Large Collections of Linguistic Resources



There is a large number and variety of linguistic resources such as corpora, lexicons, grammars, tools or experimental data. These are valuable data, but – so far – there is no established infrastructure in place to manage this data.
Often, research data is hardly accessible from outside the institutional boundaries, and sometimes, research is unnecessarily duplicated because of this.

### Complex Descriptional Requirements



There are different types of resources. While they share a common set of descriptors, they also require type-specific metadata fields.
The Dublin Core metadata set lacks descriptional means to give a full account of research data in linguistics. A more expressive metadata framework is needed such as the Component Metadata Infrastructure (CMDI)
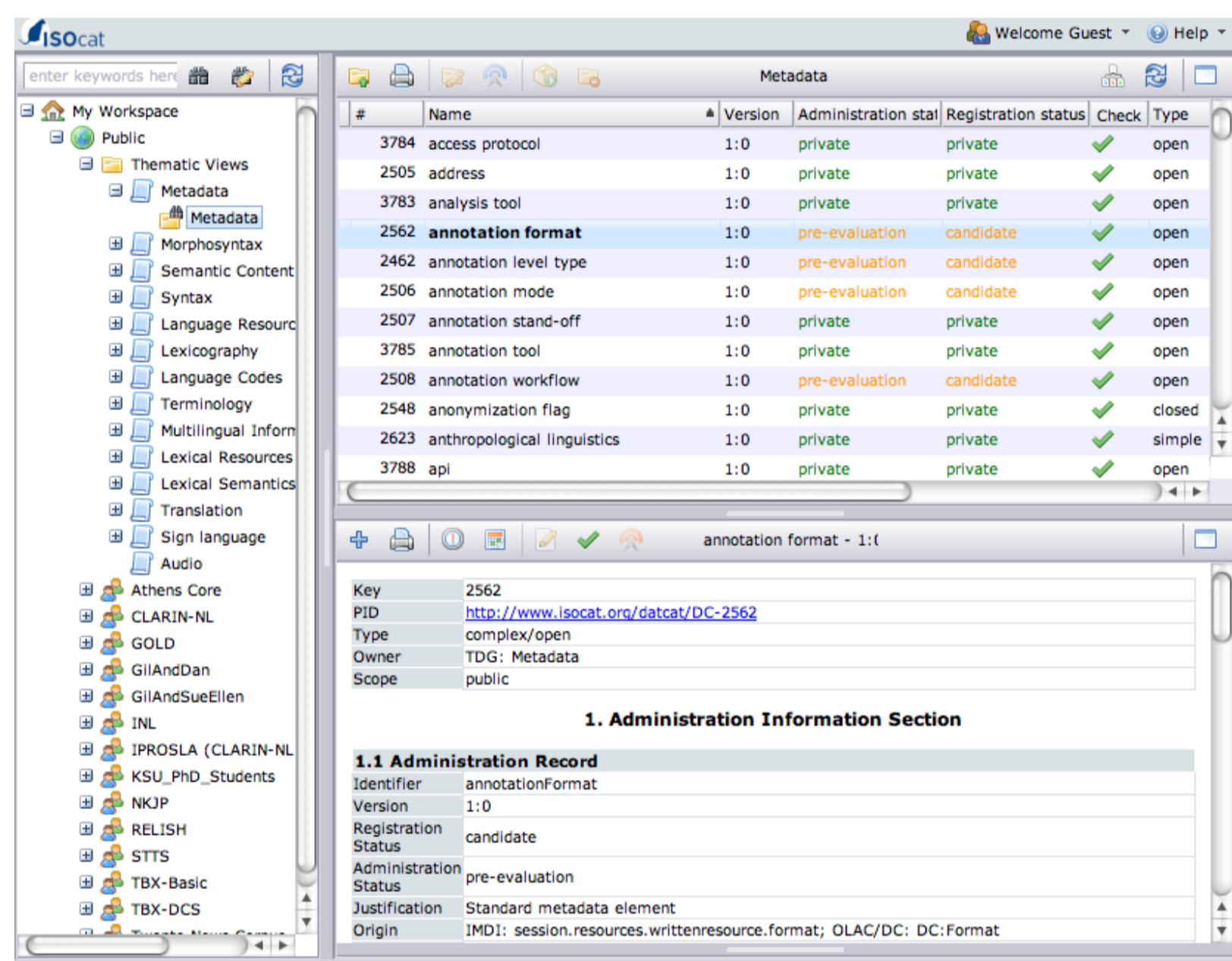
### Policy of the German Research Foundation

1. Primary research data is data that has been gathered in the course of the study of sources, experiments, measurements, and census or poll activities. They constitute the foundation for scientific publication.
2. It is necessary to devise a domain-specific organizational concept that defines the sustainable management of (research) data in a given scientific discipline.
3. The management of research data must is carried out in the framework of existing standards.
4. Data mark-up must consider rights management and include the names of data creators. An open access policy is advocated.
5. Data should become public at the end of a research project (or some limited time thereafter). This holds, in particular, for publicly funded projects. When data is attractive for commercial exploitation, divergent rules apply, but in consensus with all concerned scientists.
6. Research data must be labeled, at least, following Dublin Core. The metadata description should also include aspects relevant to the creation and subsequent reuse of research data.
7. The scientific discipline producing the research data shall devise criteria and methods to assure its quality.

Original German document at http://www.dfg.de/download/pdf/foerderung/programme/lis/ua_inf_empfehlungen_200901.pdf
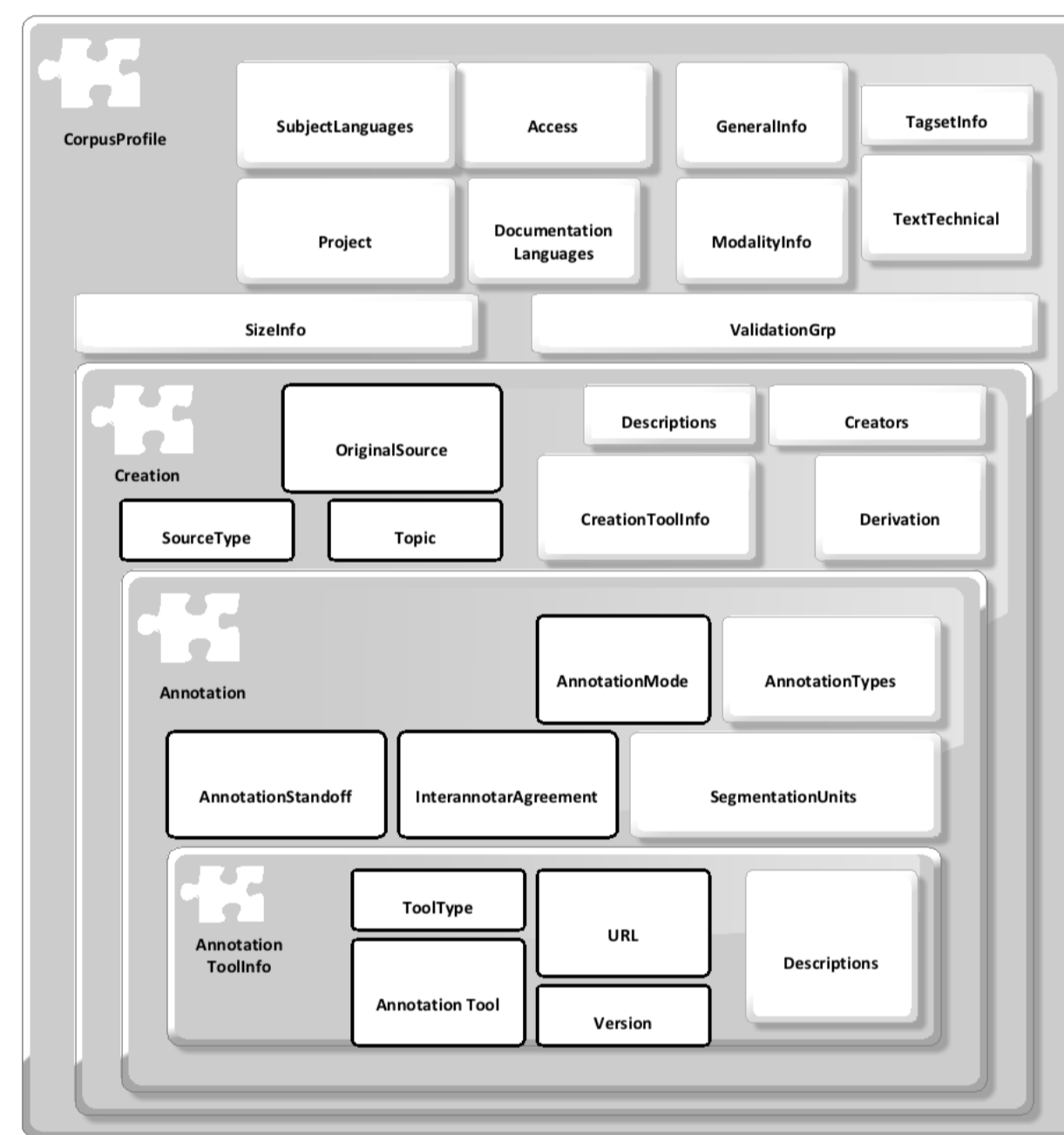
## Component Meta Data Infrastructure (CMDI)

### ISOcat Data Registry



The ISOcat registry is a community-based platform for managing elementary field descriptors for research data in linguistics. The registry has a rather flat structure, with data categories being grouped into *Thematic Domains* (e.g., *Syntax*, *Semantics*, *Terminology*, *Metadata*). There are, however, data categories of type *complex* being definable in terms of data categories of type *simple*. The thematic domain *Metadata* lists about 450 data categories.
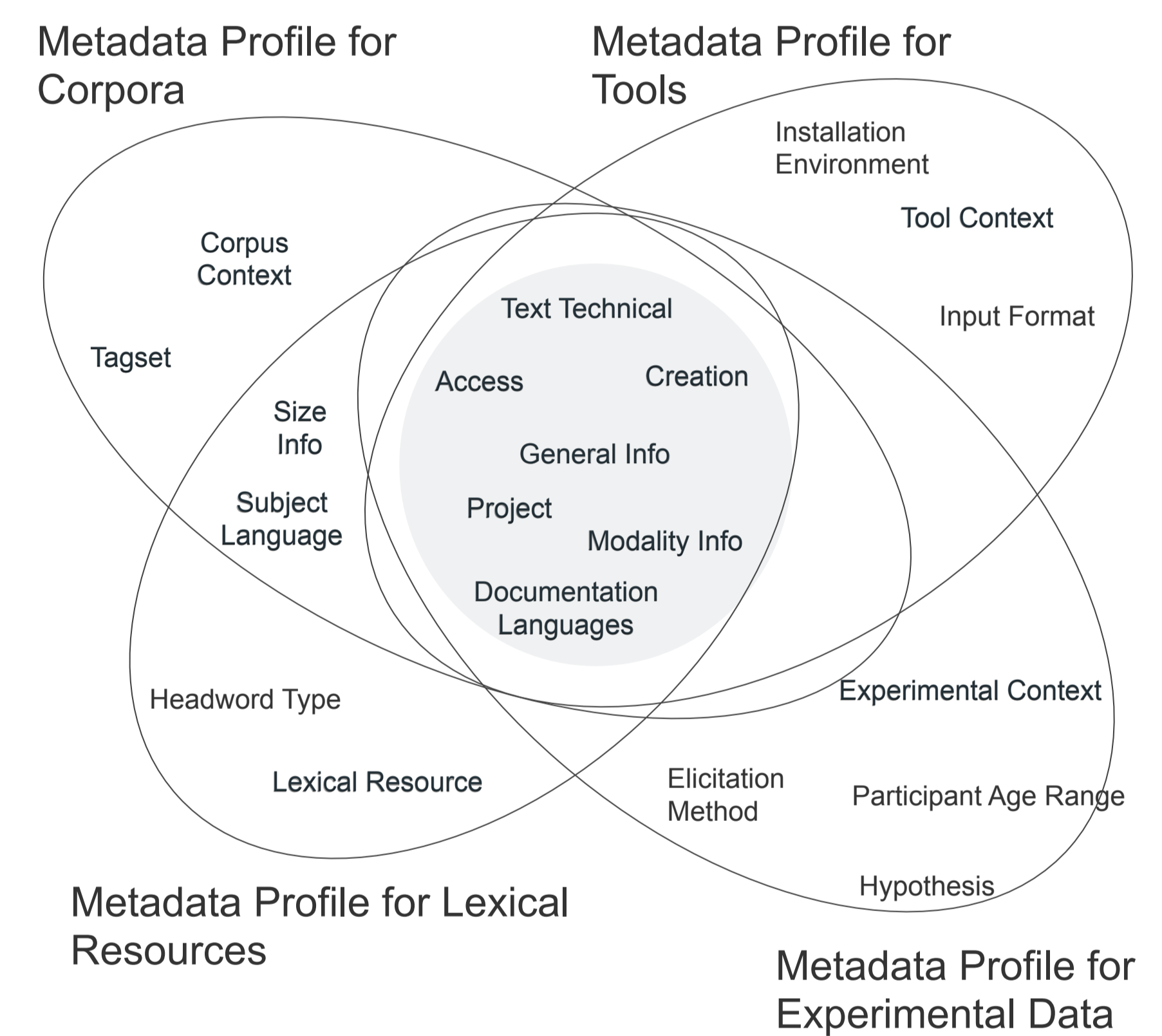Access the ISOcat registry via http://www.isocat.org.

### Component Registry



The Component Registry is a community-based platform for managing metadata structures, grouping together elementary field descriptors from the ISOcat registry or other (predefined) components. Components serve as complex building blocks for *profiles*. A profile is a schema for the description of a single type of linguistic resource.

Access the Component Registry via http://catalog.clarin.eu/ds/ComponentRegistry/#

### Profiles



In NaLiDa/SFB833, we have defined 6 profiles, and contributed many entries to the ISOcat registry and the Component Registry.

## Use of Linked Data

### Authority Files of the German National Library

The Gemeinsame Normdatei (GND) dataset pools *three German authority files:*
• Personennamendatei (**Personal Name Authority File**, PND)
• *Gemeinsame Körperschaftsdatei* (**Corporate Body Authority File**, GKD)
• *Schlagwortdatei* (**Subject Headings Authority File**, SWD)

together into a single universal authority file (GND) with > 40 Mio. triples, > 40k links to dbpedia, 38k links to lcsh, and approx. 1.8 mio links to viaf.



• **GKD**: 915.000 records of institutions (2300 entries for Tübingen institutions)
• **PND**: 3.600000 records for persons (1.800000 individualised entries)
• **SWD**: 600.000 descriptors and 700.000 synonyms, 115.000 hierarchical and 26.000 associative relations; terms grouped into 500 classes (36 clusters); 40 different subfields for linguistics; mapping of SWD entries to Dewey classification

### Metadata Editor
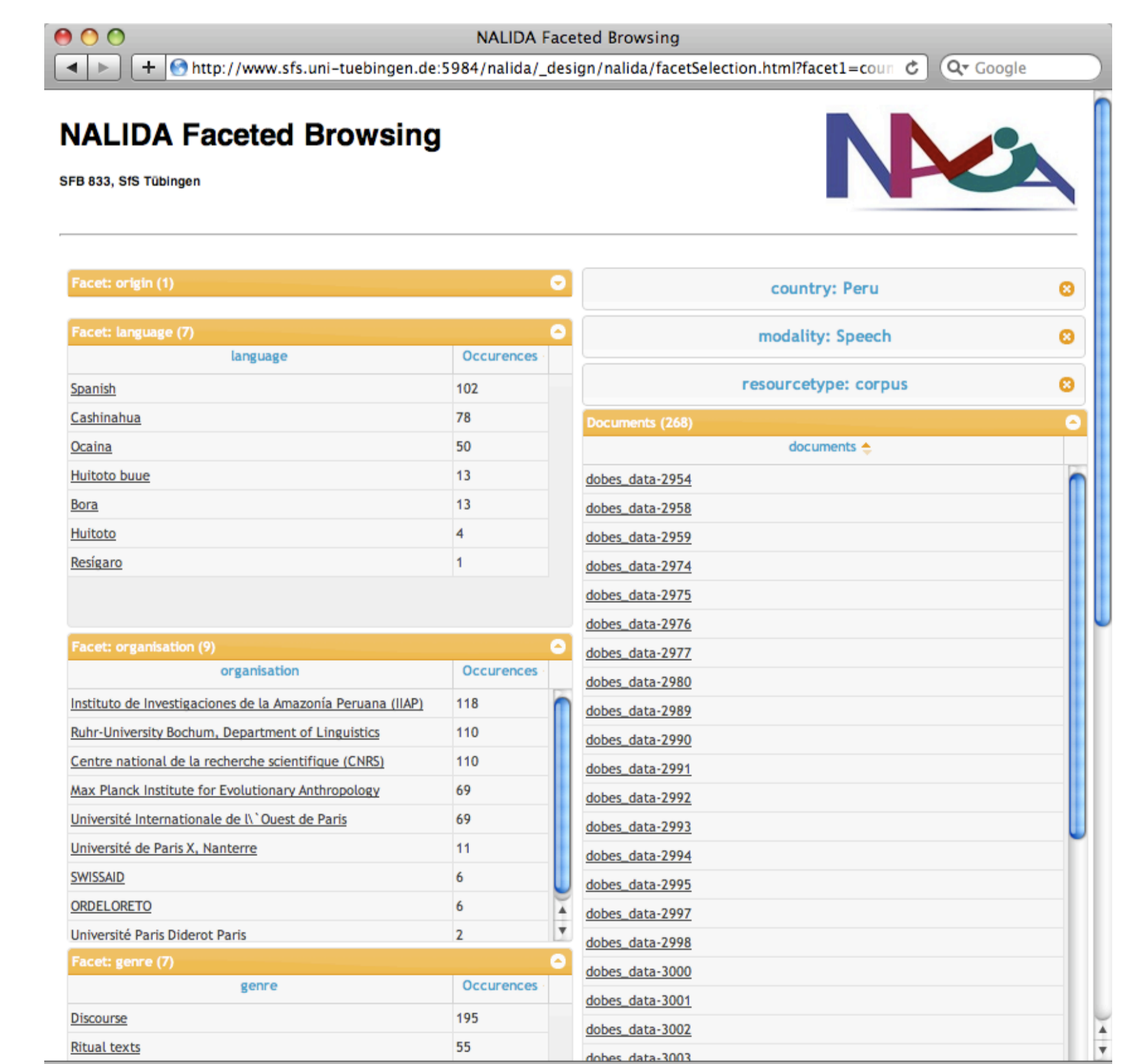


The Metadata editor is based on XForms. All XForms specifications are automatically generated from a resource's CMDI profile and a configuration file (using XSLT).
The Metadata editor offers GUI elements where values can be entered using auto-completion on PND and GKD data. Any value selected will get associated to the persistent URL of the German National Library.

**See http://www.sfs.uni-tuebingen.de/nalida/en/ for the catalogue (faceted search) and more information about the NaLiDa project.**

### Faceted Search



For the NaLiDa faceted browser, we have curated existing metadata collections in a semi-automatically manner using the GND triples.

## Acknowledgements