EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

**Faculty of Humanities**
**Department of Linguistics (SfS)**
**General and ompCutational Linguistics**

Presented at:
*Balisage* 2011
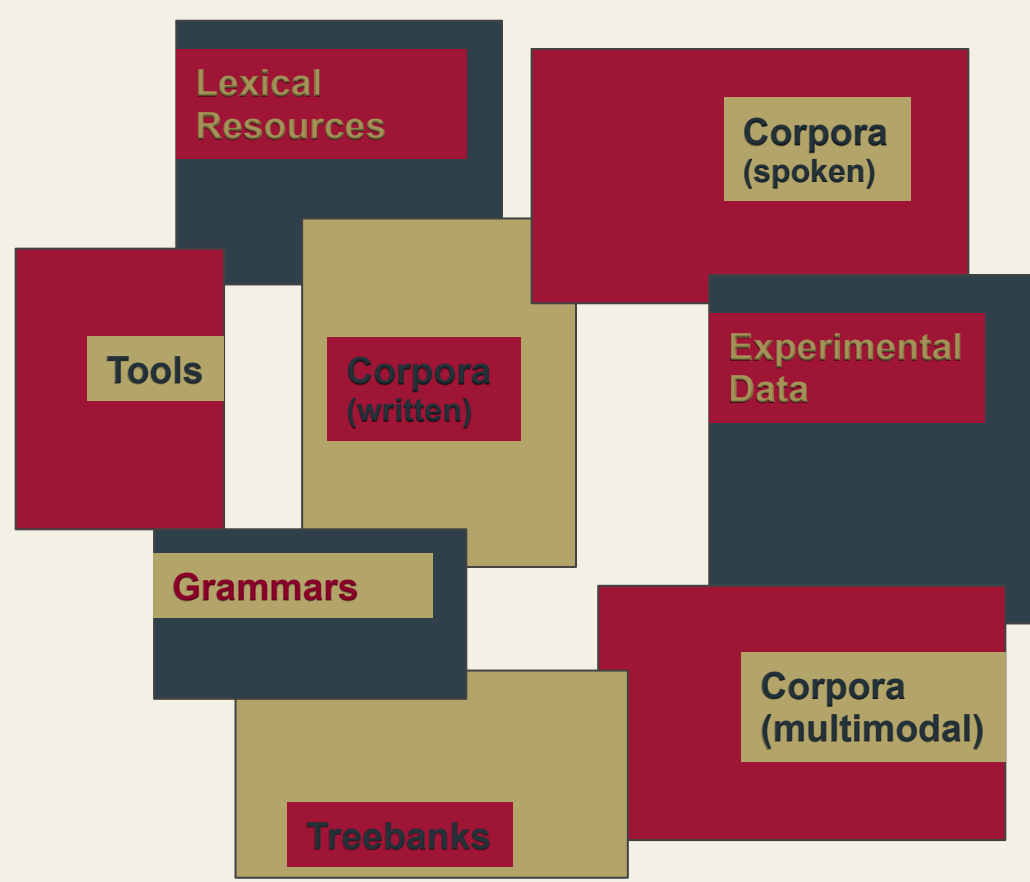The Markup Conference
August 2011, Montréal, Canada

# SUSTAINABILITY OF LINGUISTIC RESOURCES

Erhard Hinrichs, Christina Hoppermann, Thorsten Trippel, Claus Zinn

*"A language resource is sustainable if it is findable and usable independent of the development of technical standards, metadata schemes and contact persons creating or maintaining a resource."*
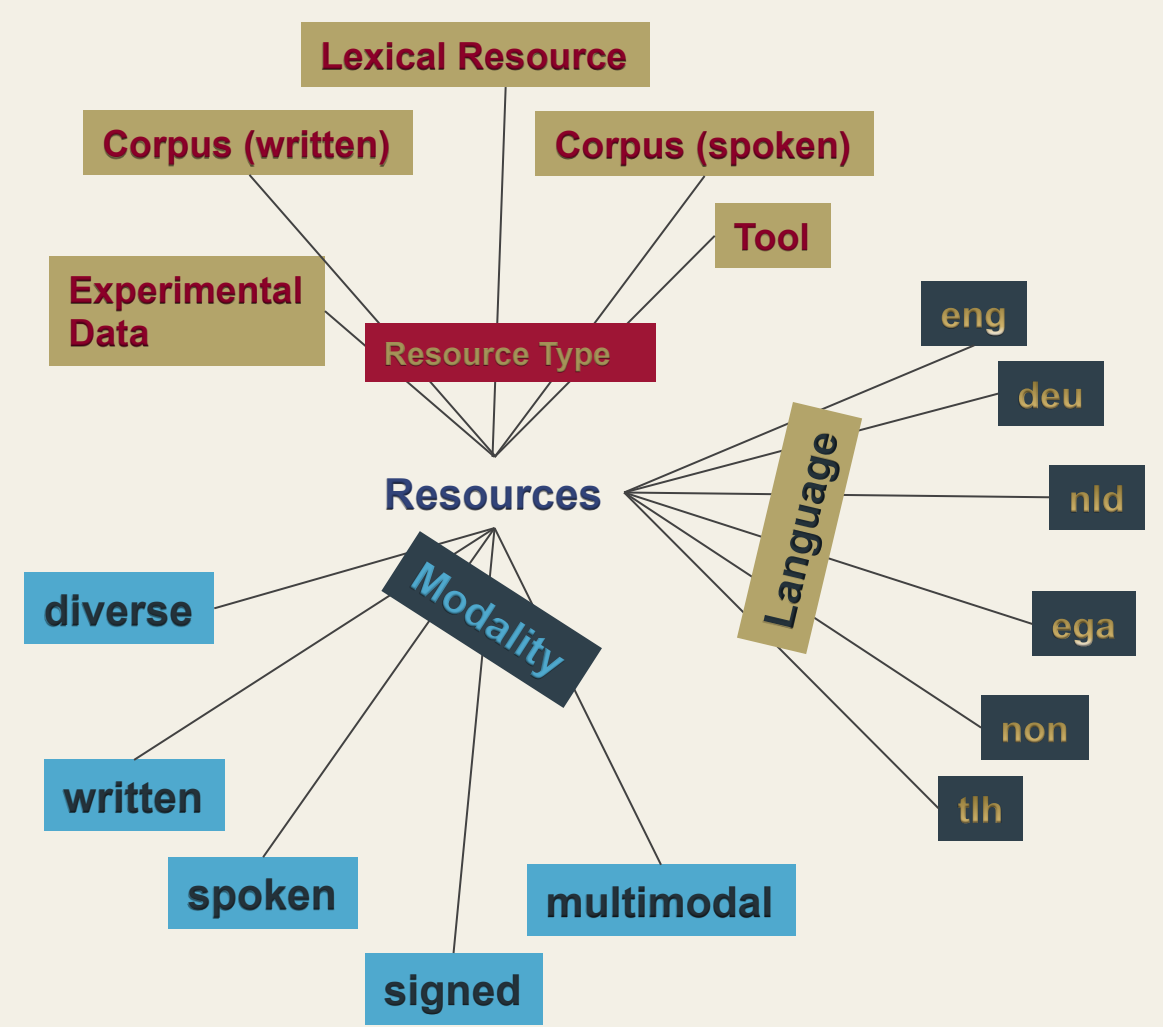
*„Eine Ressource ist nachhaltig, wenn sie auch in Zukunft trotz der möglichen Veränderungen technischer Standards, Metadatenschemas oder der nicht mehr präsenten Kontaktpersonen, die eine Ressource erstellt oder verwaltet haben, noch auffindbar und verwendbar ist."*

(From: Glossary of language resources, NaLiDa-Projekt)

For the purpose of archiving and searching, metadata are available in structured XML formats. These XML formats are also dependent on the class of a resource. A few levels of description are kept as general as possible and can thus be applied to a number of linguistic resources. Nevertheless, these general descriptions also surpass traditional bibliographies, as bibliographic metadata such as Dublin Core are of limited use for linguistic resources.

http://www.sfs.uni-tuebingen.de/nalida

Linguistic resources such as corpora, lexicons, grammars, tools or experimental data become more and more important in the linguistic research community. Their creation is often complex, bits of information get lost or cannot be processed anymore. To avoid these circumstances, resources are archived in repositories. Thereby, they are made both searchable and (re-)usable on the long run.

As there are different types of resources, there is also a need of different approaches to describe these resources. A corpus requires other descriptive elements than a tool, for instance. The assignment of a type to a resource is not necessarily unique and depends on classification criteria.

Independent of the resource type, the infrastructure needs to be able to handle and store these data types for the long term. For this purpose, there are sustainability platforms that can be used for data management, user access and as an interface to search engines. The users then only need to save their data to the platform.

The classification of resources is reflected by the search options that are based on the structures of the resource descriptions, i.e. the metadata. By this, a search for resources is restricted by characteristics such as modality, language or genre, reducing the size of the search space. Depending on selections, new structures may also appear. Here for example, *annotation schemes* were introduced based on the selection of the resource type *corpora*, a structure not relevant for many resource types.

Technical views to metadata are only relevant for computer programs (and programmers). The needs of linguists cannot be satisfied with them. The NaLiDa project offers a web portal for language resources with linguists as its target user group. Users, who are interested in language resources or their sustainable provision, can get an overview over existing resources represented by a language catalog. This catalog is based on the structures of the metadata and bundles different perspectives and aspects of the descriptions. One way of a user-friendly visualization offers the use of tabs.
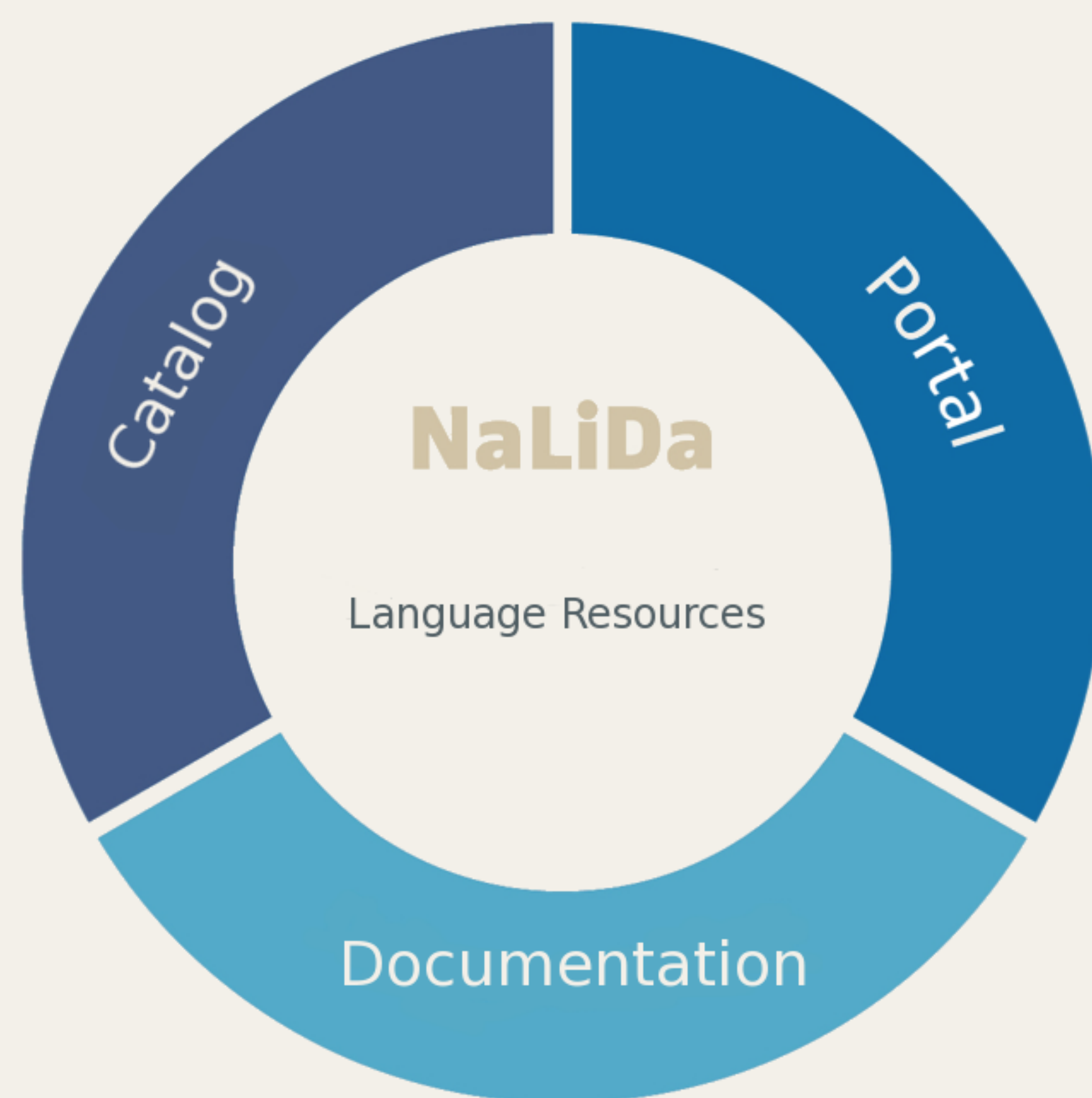
*"Research infrastructures contribute significantly to gaining knowledge in the humanities, when it comes to problems in society or accessing cultural heritage."*
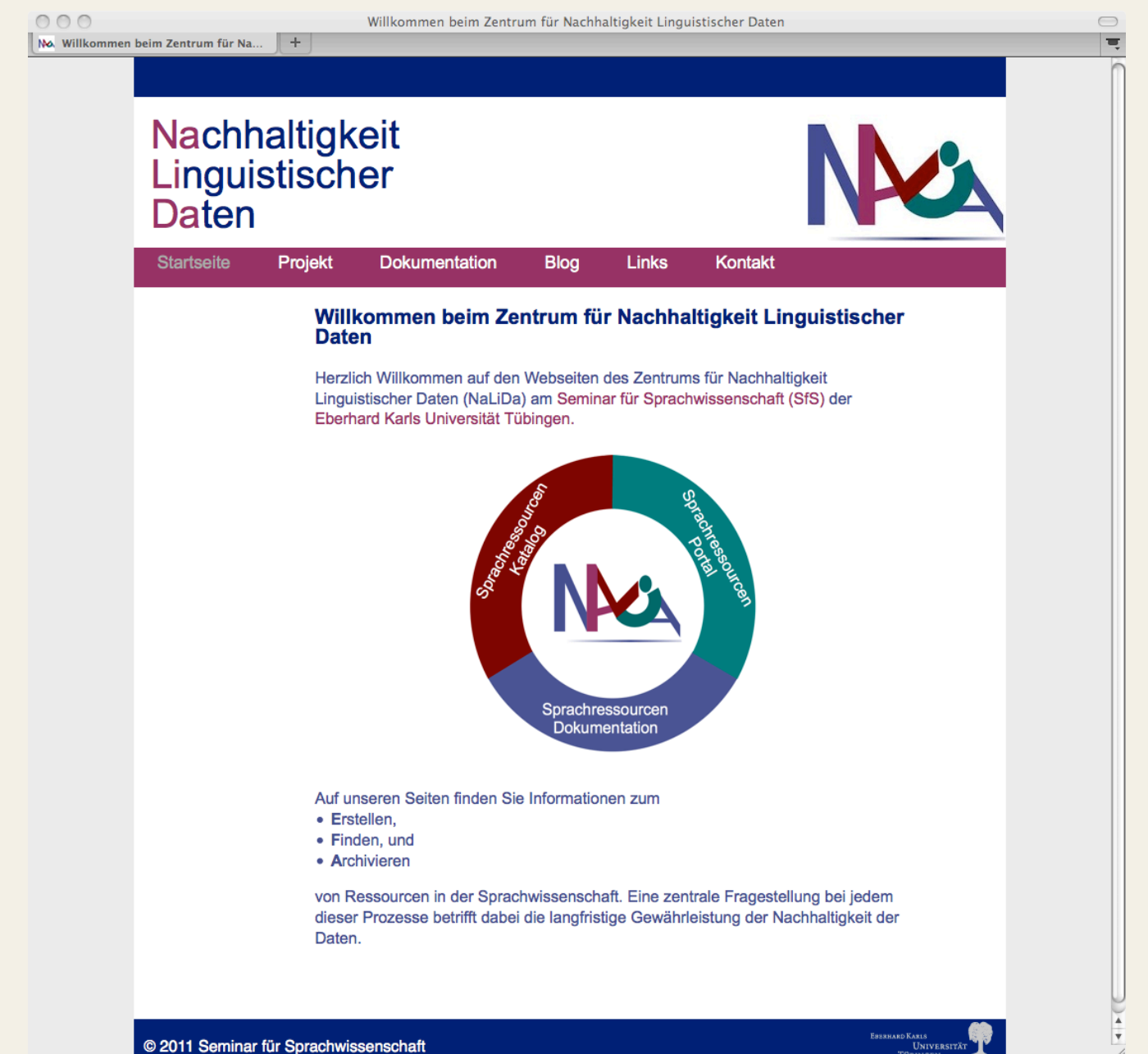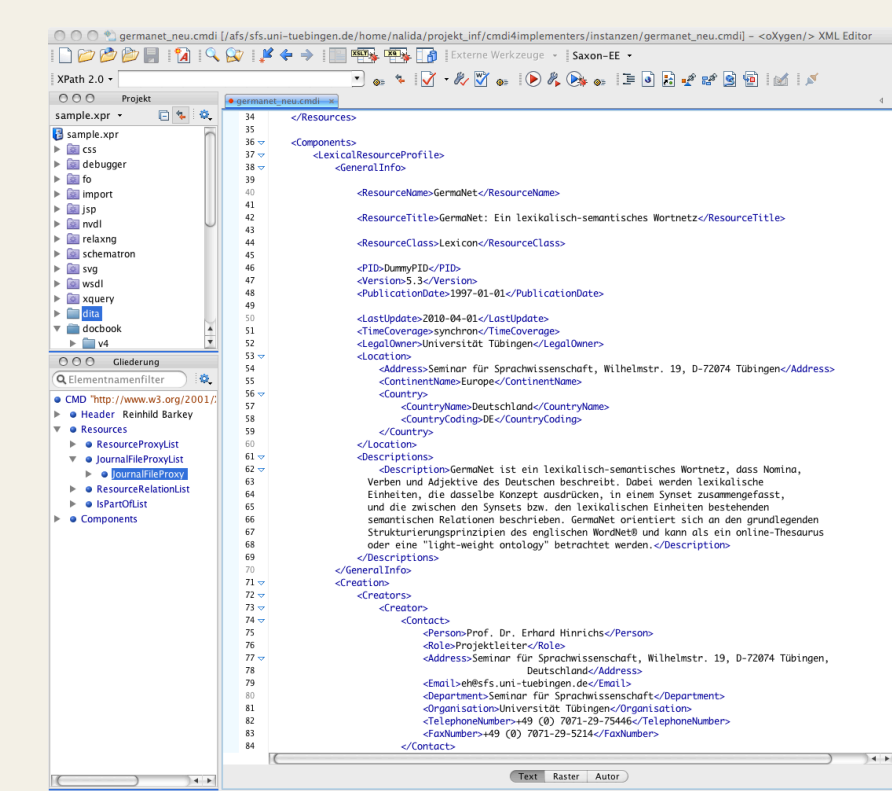
*„Forschungsinfrastrukturen leisten in den Geistes- und Sozialwissenschaften einen wichtigen Beitrag zum Erkenntnisgewinn über gesellschaftliche Problemlagen und zur Erschließung unseres kulturellen Erbes."*

(From: *Forschungsinfrastrukturen für die Geistes- und Sozialwissenschaften stärken*, Press release The German Council of Science and Humanities (Wissenschaftsrat), No. 3, Berlin, 31.01.2011)

NaLiDa — Language Resources
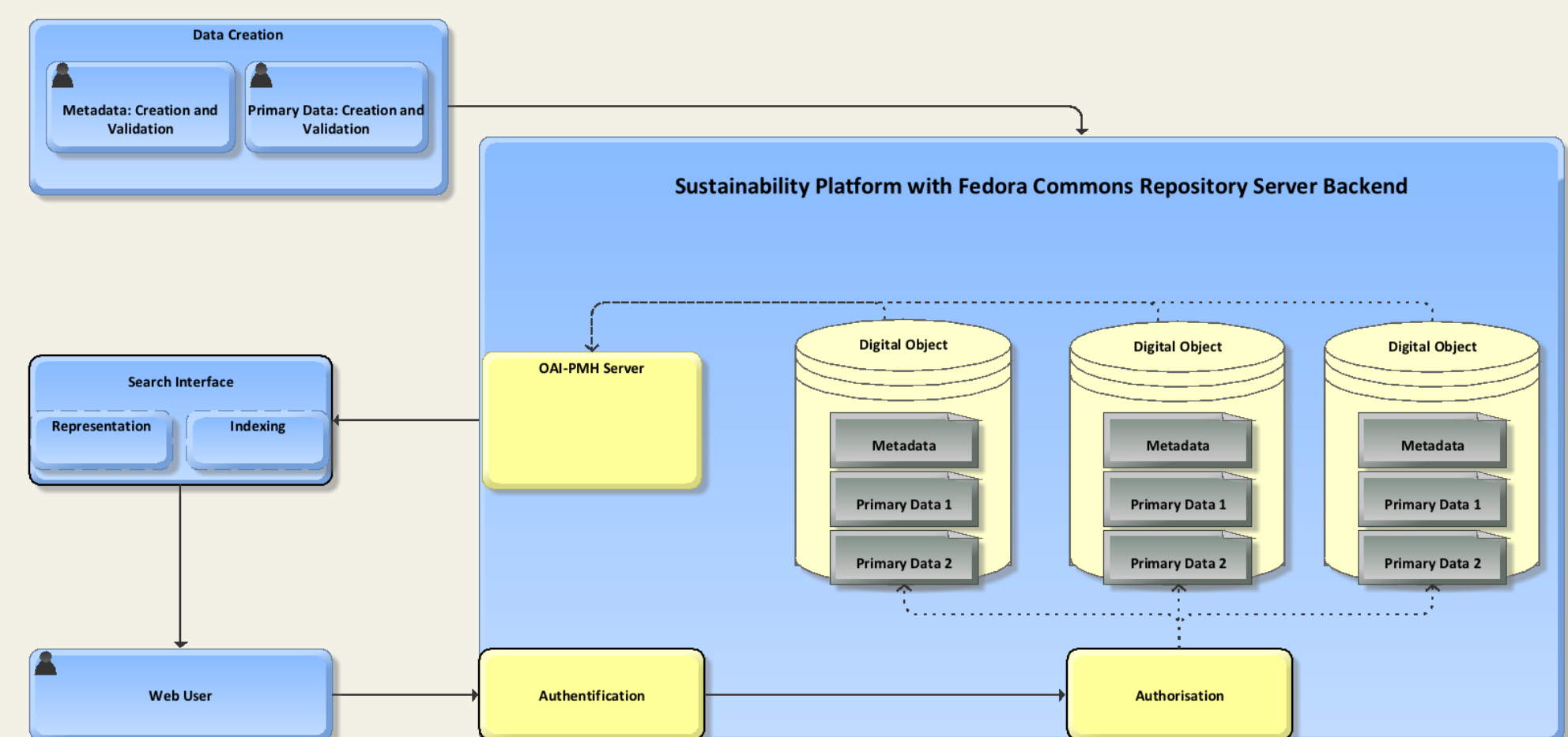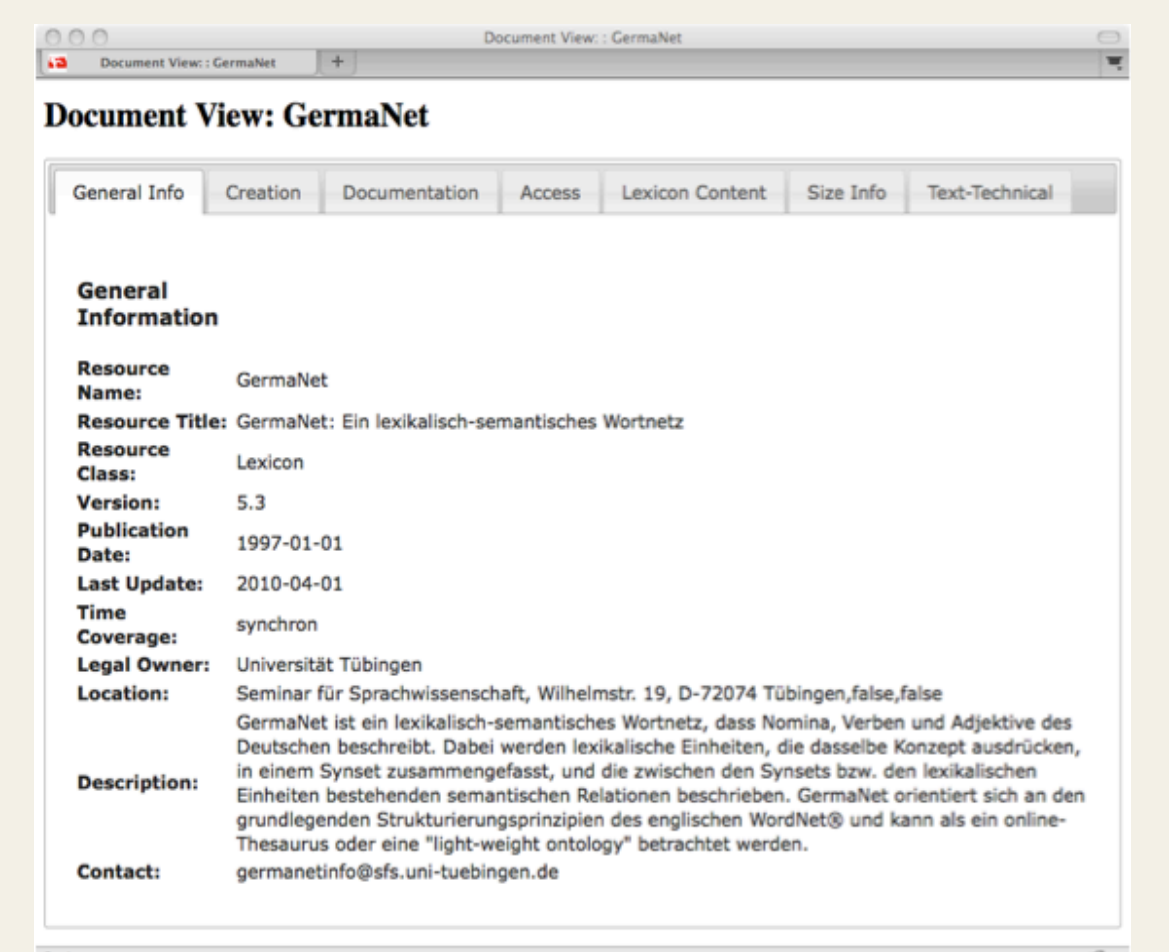(Catalog / Portal / Documentation)

*"The German Council of Science and Humanities sees the extensive public collection and provision of primary research data also as an appropriate means for quality assurance in academic practice[...]"*

*„Der Wissenschaftsrat begreift die umfassende öffentliche Sammlung und Bereitstellung von Forschungsprimärdaten auch als ein probates Mittel der Qualitätssicherung in der wissenschaftlichen Praxis [...]"*
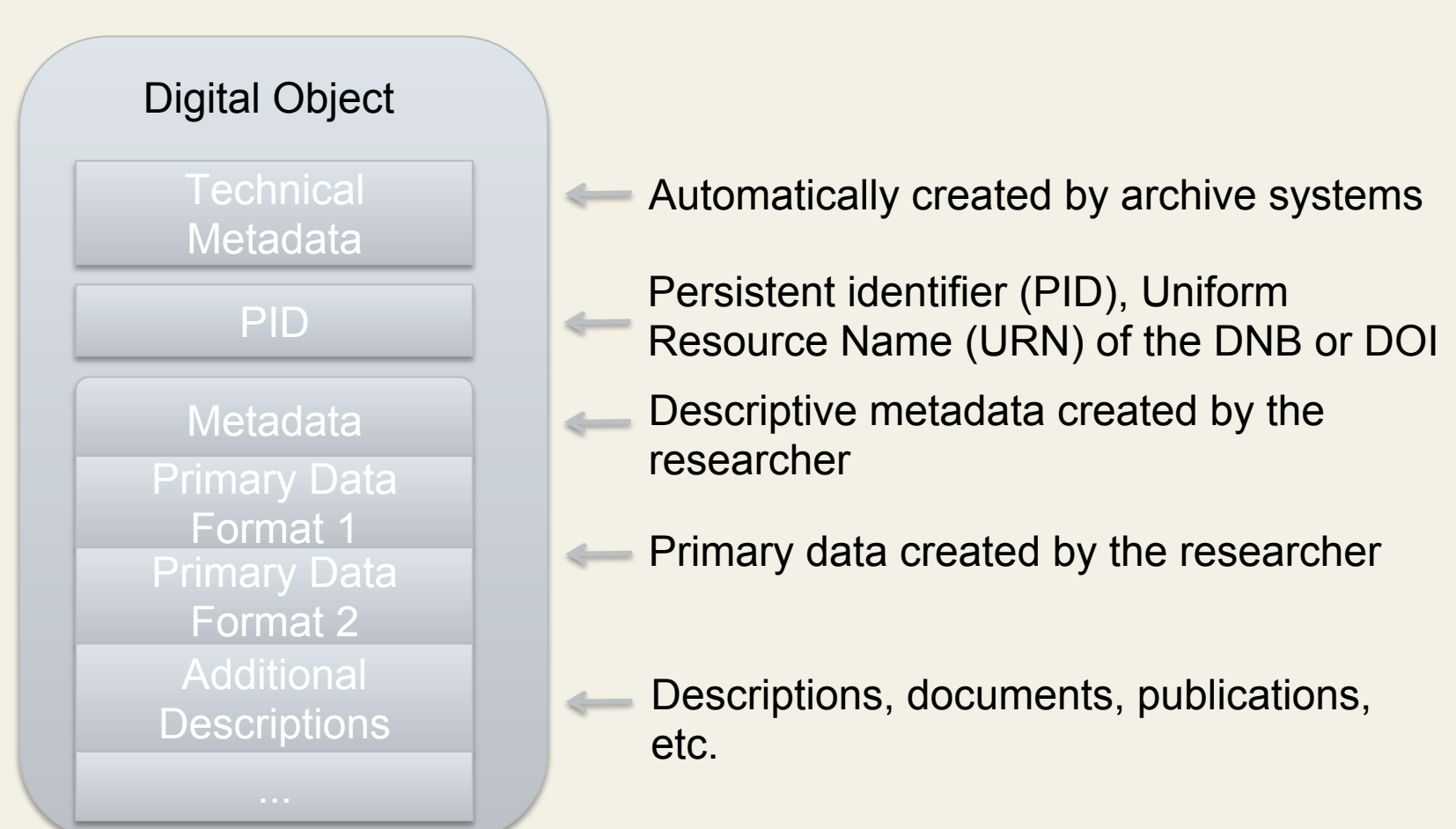
(From: Empfehlung zu Forschungsinfrastrukturen in den Geistes- und Sozialwissenschaften (Recommendations for research infrastructures in the humanities), The German Council of Science and Humanities (Wissenschaftsrat), Berlin, 28.01.2011, p.59)

*"Primary data as the basis for publications are to be kept on persistent and secured media within the institution where they were created for ten years."*

*„Primärdaten als Grundlagen für Veröffentlichungen sollen auf haltbaren und gesicherten Trägern in der Institution, wo sie entstanden sind, für zehn Jahre aufbewahrt werden."*

(From: Empfehlungen zur gesicherten Aufbewahrung und Bereitstellung digitaler Forschungsprimärdaten (Recomendations for safekeeping and providing digital primary research data), DFG: Ausschuss für Wissenschaftliche Bibliotheken und Informationssysteme, Unterausschuss für Informationsmanagement, German Research Foundation p. 2, January 2009)

The search functions and metadata views within the web portal provide access to primary research data such as corpora, lexical resources, experimental data, etc. Privacy and property restrictions may require limited access to the public, hence authorization is part of the archive system. Reasons for restrictions can include: pending publication approval, ethical reasons, etc. Nonetheless, such resources can still be found on the basis of metadata so that it is at least possible to approach the resource creators and learn about the conditions for accessing the resource.

The essential elements of the sustainability platform are digital objects. A digital object is a collection that contains a resource as well as its metadata description or any additional information belonging to it. This information may include scientific research articles that either describe the resource or use it as a basis, internal (project) reports, notes, etc. For authorisation restrictions on user groups, the platform supports access limitations.
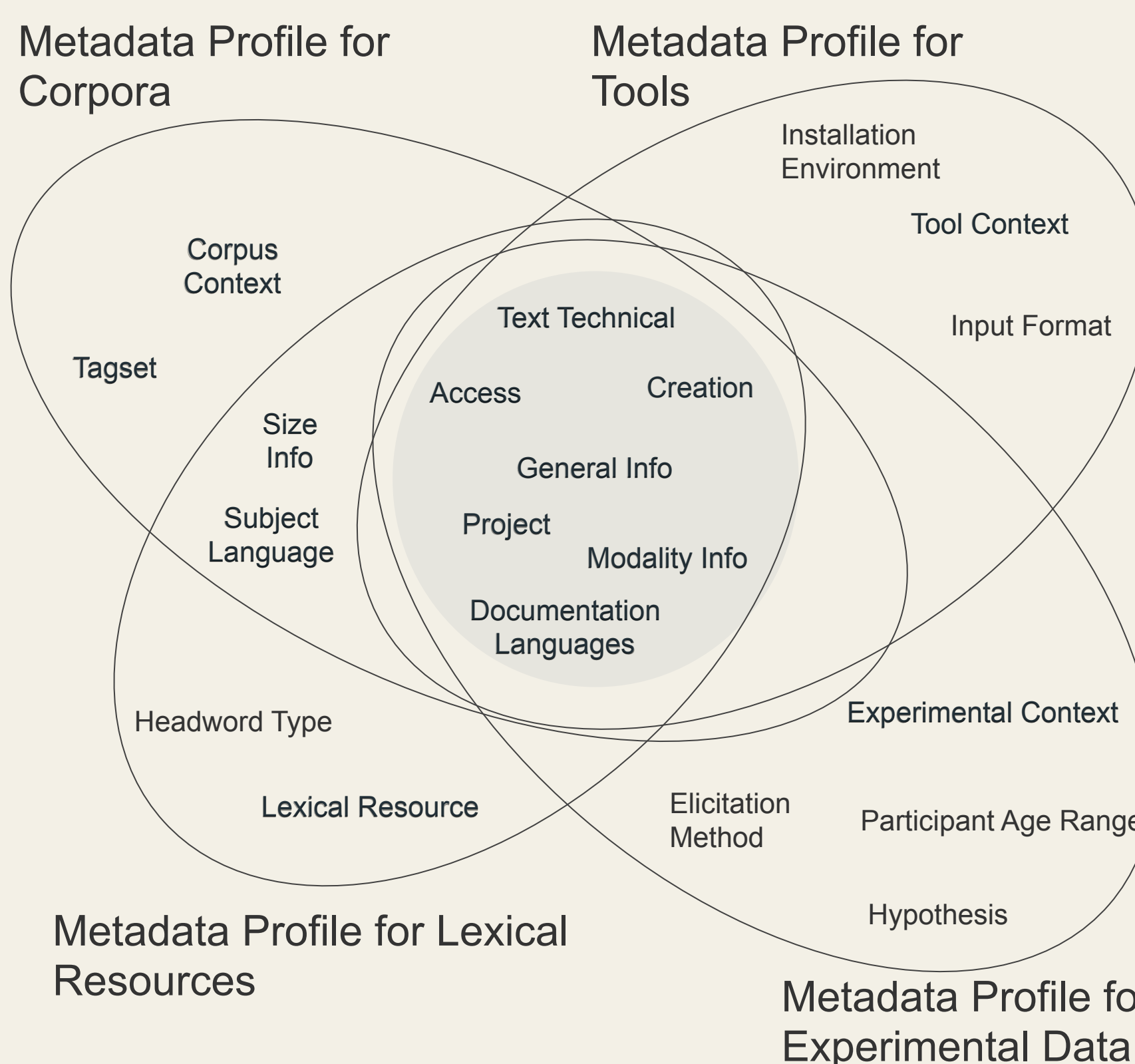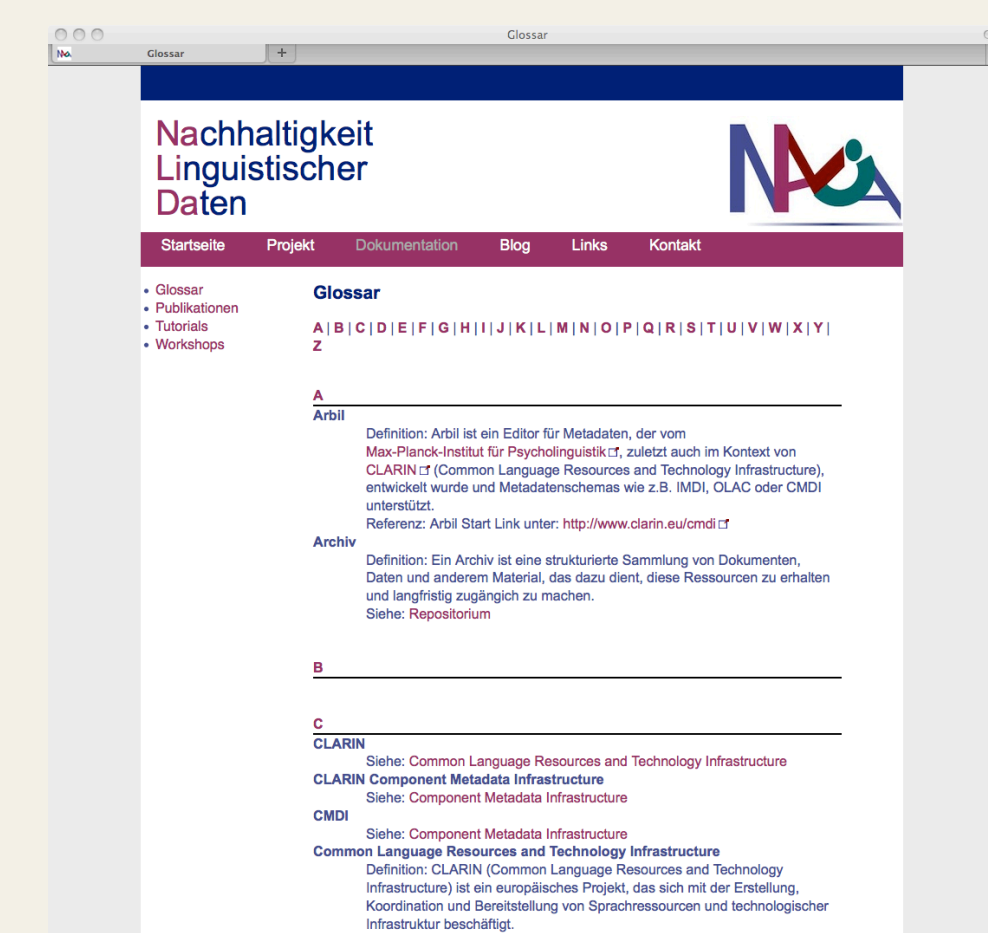
Digital objects can be identified and addressed by unique persistent identifiers (PIDs). Libraries, such as the German National Library (DNB), have developed ranges and formats for these identifiers. Other examples are DOIs.

The provision of linguistic resources and the search over these, is a major part of the web portal, also serving the purpose of exchanging information on the creation and archiving of language resources. In order to support the communication between different research branches, there is - apart from other documentations and tutorials - a glossary which makes terms accessible that are relevant for linguistic resources. This is also part of the sustainability aspect.

Within the NaLiDa project, the creation of metadata generally uses an iterative process. This mainly serves the purpose of training resource creators enabling them to create metadata on their own, but also of including new resource types. On the basis of first interviews and resource descriptions, a draft of a metadata schema is created together with a metadata instance. If needed, further components and profiles are also developed.

This metadata instance can then be corrected by the resource creator without requiring any further technical knowledge. Adjustments can be made at any time.

If the metadata instance is complete and correct, it can be stored in the sustainability platform together with the primary data. After this process, it is still possible to edit the instance.

Interview → Draft metadata → Revision → Repository

Metadata Profile for Corpora
Metadata Profile for Tools
Metadata Profile for Lexical Resources
Metadata Profile for Experimental Data

Corpus Context / Tagset / Size Info / Subject Language / Headword Type / Lexical Resource / Text Technical / Access / Creation / General Info / Project / Modality Info / Documentation Languages / Installation Environment / Tool Context / Input Format / Experimental Context / Participant Age Range / Elicitation Method / Hypothesis

Digital Object
Technical Metadata — Automatically created by archive systems
PID — Persistent identifier (PID), Uniform Resource Name (URN) of the DNB or DOI
Metadata — Descriptive metadata created by the researcher
Primary Data Format 1 / Primary Data Format 2 — Primary data created by the researcher
Additional Descriptions — Descriptions, documents, publications, etc.

The long-term archiving does not only include the technical storage, i.e. the storing of primary data in so-called data streams, but also the descriptions of the resources, i.e. the metadata. Within the NaLiDa project, CMDI (Component Meta Data Infrastructure) is used as the underlying metadata schema. For each resource class there is a metadata profile available for describing a resource. The smallest building blocks, i.e. the data categories (DatCats), are collected in components that can be reused in different profiles independent of the resource class. Often-used components build a core which is comparable to Dublin Core, but includes more data categories relevant for linguistic resources.

University of Tübingen· Faculty of Humanities· Department of Linguistics
Wilhelmstraße 19 · 72074 Tübingen · Germany
www.sfs.uni-tuebingen.de/nalida
nalida@sfs.uni-tuebingen.de