# Metadata and the use of it

# Unterrichtseinheit Metadaten im Rahmen des Austausches von Lehrmaterial

- Zielgruppe der Lehreinheit:
  - Linguisten und fortgeschrittene Linguistik-Studierende
  - Computerlinguisten und –studierende
  - Z.T. praktizierende Wissenschaftler, die Daten archivieren wollen oder müssen
- Ziele:
  - Teilnehmer verstehen, was Metadaten zur Beschreibung von Ressourcen sind
  - Teilnehmer kennen Gründe, warum Metadaten sinnvoll und nützlich sind
  - Teilnehmer können Metadatenformulare ausfüllen
  - Teilnehmer kennen Metadateninfrastrukturen: ISOcat, Komponenten, Profile
  - Teilnehmer haben eine Vorstellung, wie Metadaten verbreitet werden
- Die Folien und Auszüge daraus dürfen für die eigene Erstellung von Lehrmaterialien verwendet werden

# Aufbau dieser Unterrichtseinheit

- Der Foliensatz besteht aus zwei Versionen:
  - Foliensatz für den Dozenten: hier sind zusätzliche Informationen zum Aufbau und zur Struktur enthalten, die für Teilnehmer nicht primär hilfreich wären
  - Foliensatz für den Unterricht: die Folien dienen der sequentiellen Präsentation. Auf Effekte wird verzichtet, so dass die Präsentation auch als PDF erfolgen kann.
- Um technische Herausforderungen zu vermeiden, ist als Zielformat PDF vorgesehen. Dadurch werden Präsentationsfunktionen wie Referentenansicht, etc. nicht verwendet.
- Die Unterrichtseinheit ist in "Kapitel" untergliedert, am Ende jedes Kapitel könnten Pausen eingeschoben oder die Sitzung beendet werden. Kapitel werden als Einheit betrachtet, sie sind aber unterschiedlich lang.
- Übungsaufgaben werden angeboten, es wird empfohlen sich zumindest eine Lösung zu überlegen. Die Lösungen werden nicht immer eindeutig sein, im Dozentenfoliensatz ist immer eine Beispiellösung enthalten.
- Informationen aus Tutorien vom Projekt NaLiDa (http://www.sfs.uni-tuebingen.de/nalida/de/doku/tutorials) und CLARIN (www.clarin.eu/cmdi) werden dabei verwendet. Daher geht ein besonderer Dank für die extensive Dokumentation an Christina Hoppermann, Daan Broeder, Dieter van Uytvanck, Menzo Windhouwer.

- Dublin Core

- Bibliographie

- Datenkategorie

- DatKat

- ISOcat

- Datenkategorien-Repositorium

- Repositorium

- OAI

- OLAC

- Component Metadata Infrastructure

- CMDI

- Komponenten

- Profile

- Instanzen

- Ressource

# Bibliographie (Auswahl)

- CMDI, http://www.clarin.eu/cmdi
- CMDI Broeder, D., Schonefeld, O., Trippel, T., Van Uytvanck, D., and Witt, A. (2011). A pragmatic approach to XML interoperability — the component metadata infrastructure (CMDI). In *Balisage: The Markup Conference 2011*, volume 7.
- I. Schuurman, M.A. Windhouwer. Explicit Semantics for Enriched Documents. What Do ISOcat, RELcat and SCHEMAcat Have To Offer? 2nd Supporting Digital Humanities conference (SDH 2011), 17-18 November 2011, Copenhagen, Denmark.
- Van Uytvanck, D., Zinn, C., Broeder, D., Wittenburg, P., & Gardelini, M. (2010). Virtual language observatory: The portal to the language resources and technology universe. *Proceedings of the Seventh conference on International Language Resources and Evaluation (LREC'10)* (pp. 900-903). European Language Resources Association (ELRA).
- *Gavrilidou M., Labropoulou P., Piperidis S., Monachini M., Frontini F., Francopoulo G., Arranz V., and Mapelli V. (2011). A* Metadata Schema for the Description of Language Resources (LRs), 5th International Joint Conference on Natural Language Processing (IJCNLP 2011).

- D. Broeder, D. van Uytvanck, M. Gavrilidou, and T. Trippel. Standardizing a component metadata infrastructure. In Proceedings of the 8th Conference on International Language Resources and Evaluation (LREC 2012), Istanbul, Forthcoming.
- CMDI MI Search Engine. CMDI Meertens Institute Search Engine. http://www.meertens.knaw.nl/cmdi.
- M. Durco, D. Broeder, and M. Windhouwer. Semantic mapping - groundwork for query expansion and semantic search. 2012, submitted.
- ISO 12620:2009. Terminology and other language and content resources - specification of data categories and manage- ment of a data category registry for language resource. ISO, 2009.
- ISOcat. http://www.isocat.org.
- OLAC. http://www.language-archives.org/.
- VLO. http://catalog.clarin.eu/ds/vlo/.
- Dublin Core: http://dublincore.org/
- Tutorials (z. B. Component Registry): http://www.sfs.uni-tuebingen.de/nalida/de/doku/tutorials

# Struktur

- Einleitung
- Motivation
- Beschreibung eines Buches: Bibliographien und Dublin Core
- Beispiele für linguistische Ressourcen: Beschreibung durch OLAC
- Datenkategorien als Beschreibungsebenen für Ressourcen
- Zentrales Verzeichnis für Datenkategorien: ISOcat und Beschreibung einer Ressource mit ISOcat Datenkategorien
- Generalisierungen von Gruppen von Datenkategorien: Komponenten
- Generalisierungen von Gruppen von Komponenten für Ressourcentypen
- Die CMDI Component Registry
- Erstellen von CMDI-Instanzen
- Zusammenfassung

# Chapter 1

Introduction of this unit

# Goal of this unit

- Introducing metadata

- Creating a metadata schema

- Adjusting metadata schemas

- Describe different types of resources

- Create metadata with

  – Data categories from the CLARIN Concept Registry

  – Components from the Component Registry (including creating and using them)

- Use tools for creating metadata

# Chapter 2

Motivation for metadata

# Metadata – what is that?

"Metadata is often described as data about data or information about information."

http://www.niso.org/publications/press/UnderstandingMetadata.pdf
(Metadata is often called data about data or information about information)

"Metadata are structured data used to describe informative resources. With this description the data can be found more easily."
http://www2.sub.uni-goettingen.de/intrometa.html

Metadata are machi[n] readable bits of information for electronic resources or other things. M
(Übersetzung siehe: http://www2.sub.uni-goettingen.de/intrometa.html, Original: http://www.w3c.org/metadata)

# Metadata – what is that?

"Metadata is often
described as data

Metadata are machine
readable

Finding (a resource):
The process to locate a
resource in a collection
of resources, using
attributes and
characteristics of the
resource for
identification
purposes., for example
the position in a
classification system.

Description (of a resource):
provides an overview for a
user, for example on
bibliographic information,
restricted usage, contact
infomration, type of resource,
etc.

"Metadata are structured
data used to describe
informative resources. With
this description the data can
be found more easily."
http://www2.sub.uni-
goettingen.de/intrometa.html

# Data and metadata in archives and repositories

- Archive:
  - Authoritative reference copy of data
  - Long(er) time storage
- Data safety and integrity by access control
  - Identification of users by the means of central infrastructures
  - Access control lists for repository access using defined procedures
  - Defined fall back procedures for discontinued projects or organizational units
  - Availability of data
- Citability by reference copy
- Citation by referencing the description of the resource Description of research data: part of the archiving workflow

- Name some archives and repositories!

- What kind of resources are "saved" there?

- What could be the metadata there?

CLARIN-D

- Zählen Sie Archive und Repositorien auf, die Sie kennen!
  - Mögliche Antworten:
  - Universitätsbibliothek
  - Stadtarchiv
  - Museum
- Was für Ressourcen werden dort "gespeichert"?
  - Bücher
  - Museale Artikel
- Was könnte man dort unter Metadaten verstehen?
  - Informationen im Bibliothekskatalog
  - Museums-Bestandsliste/Inventar

# Chapter 3

Describing a book

# What is that?

# What is that?

- Describe this book, keep in mind the aspect of
  - Locating it
  - Describing the content
- Look at your description
  - Anything familiar?
  - What is it?

nestor Handbuch
Eine kleine Enzyklopädie
der digitalen Langzeitarchivierung
Version 2.0

H. Neuroth, A. Oßwald, R. Scheffel, S. Strathmann, M. Jehn [Hrsg.]

nestor vwh

# Was ist denn das?

- Beschreiben Sie dieses Buch unter dem Aspekt der
  - Auffindbarkeit
  - Des Inhalts

  Antwort: Die Bibliographie dieses Buches:

  Herausgeber: H. Neuroth, A. Oßwald, R.

  Scheffel, S. Strathmann,  M. Jehn,

  Titel=   nestor Handbuch: Eine kleine

  Enzyklopädie der

            digitalen Langzeitarchivierung

  Jahr =   2010,

  Identifikator= urn:nbn:de:0008-20100305186

  Farbe= blau, ...

  Schlagwörter= Archivierung, digitale Ressourcen,

  ...

- Betrachten Sie Ihre Beschreibung
  - Kommt Ihnen etwas bekannt vor?
  - Was?
  - Antwort: Das ist das, was man in einem
    Bibliothekskatalog finden würde.

nestor Handbuch: Eine kleine Enzyklopädie der digitalen Langzeitarchivierung
hg. v. H. Neuroth, A. Oßwald, R. Scheffel, S. Strathmann, M. Jehn
im Rahmen des Projektes: nestor – Kompetenznetzwerk Langzeitarchivierung und
Langzeitverfügbarkeit digitaler Ressourcen für Deutschland
nestor – Network of Expertise in Long-Term Storage of Digital Resources
http://www.langzeitarchivierung.de/

Kontakt: editors@langzeitarchivierung.de
c/o Niedersächsische Staats- und Universitätsbibliothek Göttingen,
Dr. Heike Neuroth, Forschung und Entwicklung, Papendiek 14, 37073 Göttingen

Die Herausgeber danken Anke Herr (Korrektur), Martina Kerzel (Bildbearbeitung) und
Jörn Tietgen (Layout und Formatierung des Gesamttextes) für ihre unverzichtbare
Unterstützung bei der Fertigstellung des Handbuchs.

Bibliografische Information der Deutschen Nationalbibliothek
Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen
Nationalbibliografie; detaillierte bibliografische Daten sind im Internet unter
http://www.d-nb.de/ abrufbar.

Die Inhalte dieses Buchs stehen auch als Onlineversion
(http://nestor.sub.uni-goettingen.de/handbuch/)
sowie über den Göttinger Universitätskatalog (http://www.sub.uni-goettingen.de) zur
Verfügung.
Die digitale Version 2.0 steht unter folgender Creative-Commons-Lizenz:
„Attribution-Noncommercial-Share Alike 3.0 Unported"
http://creativecommons.org/licenses/by-nc-sa/3.0/

Einfache Nutzungsrechte liegen beim Verlag Werner Hülsbusch, Boizenburg.
© Verlag Werner Hülsbusch, Boizenburg, 2009
www.vwh-verlag.de
In Kooperation mit dem Universitätsverlag Göttingen

Markenerklärung: Die in diesem Werk wiedergegebenen Gebrauchsnamen, Handelsnamen,
Warenzeichen usw. können auch ohne besondere Kennzeichnung geschützte Marken sein und
als solche den gesetzlichen Bestimmungen unterliegen.

Druck und Bindung: Kunsthaus Schwanheide

Printed in Germany – Als Typoskript gedruckt –

ISBN: 978-3-940317-48-3

URL für nestor Handbuch (Version 2.0): urn:nbn:de:0008-2009073109
http://nbn-resolving.de/urn/resolver.pl?urn=urn:nbn:de:0008-2009073109

- "external"
  - Title
  - Author
    - Authority file (e.g. VIAF)
  - publisher
  - year
  - Page number
  - cover
  - institution
    - Authority file again (GKD)

- "internal"
  - topic
  - keyword
    - Standard keyword files (SND)
  - Position in classification system
    - Dewey Decimal Classification (DDC)
    - Universal Decimal Classification (UDC)

# Purpose of bibliographic data: locating

- Signature or location information
- Together with a library map

# Chapter 4

Example of linguistic resources

# Examples of linguistic resources

Lexical Resources

Corpora (spoken)

Software Tools

Corpora (written)

Experimental-Data

Grammars

Multimodal Corpora

Treebanks

# Resource type: text corpus

Corpora
(written)

| | |
|---|---|
| Resource Name: | TüBa-D/Z |
| Resource Title: | Tübinger Baumbank des Deutschen/Zeitungskorpus |
| Resource Class: | Corpus |
| Corpus Type: | general corpus |
| Annotation Type: | POS, inflectional morphology, named entities, referential relations, grammatical functions |
| Annotation Mode: | semi-automatisch |
| Tagset: | Stuttgart-Tübingen Tagset (STTS) |

# Resource type: Tool

**Software Tools**

| | |
|---|---|
| Resource Name: | Tree Tagger |
| Resource Title: | Tree Tagger: A Language - Independent POS Tagger |
| Resource Class: | Tool |
| Tool Type: | Annotation Tool, Lemmatizer, Chunker, POS Tagger |
| Input: | wortsegmentierter Text |
| Output: | Annotierter Text mit Informationen zu POS und Lemmata |
| Operating System: | Linux, Windows, Mac OS X, Solaris |

# Resource type: Experimental data

|  |  |
|---|---|
| Resource Name: | Experiment 1 |
| Resource Title: | Experiment on bridge verbs and V2 structures |
| Resource Class: | Experiment |
| Experiment Type: | Survey Data |
| Domain: | Linguistik |
| Elicitation Method: | Fragebogen |
| Elicitation Software: | WebExp |
| Sample Size: | 23 Teilnehmer |
| Mean Age: | 30,4 Jahre |

Experimental-Data

- What kinds of resources can you think of?

- What types of resources do you work with?

- What type of resources do you create?

- How would you describe them?

- Welche Ressourcen kennen Sie?

- Mit welchen Ressourcen arbeiten Sie?

- Was für Ressourcen erstellen Sie selber?

- Wie würden Sie die Ressourcen beschreiben?


- Beispiele finden sich auf den vorherigen Folien

# Chapter 5

Data categories for describing resources

# Attribute-Value-Structures

- Pairs of "Name" und "Value"
    - Name: Chomsky
    - First name: Noam
    - Name: Halle
    - First name: Morris
    - Title: The sound pattern of English
    - Year:1968

- Other terminology:
    - Data category value pair
    - Name-Value-Pair
    - Key-Value-Pair
    - AV-Structures
    - ...

Name    Value

- Pairs of "Name" and "Value"
  - Name: Chomsky
  - First name: Noam
  - Name: Halle
  - First name: Morris
  - Title: The sound pattern of English
  - Year:1968

- Who are the authors?
  - Hint:
    - Ambiguous structure
    - "Order" is not meaningful in AV-Structures
  - There are two options
    - If everybody has one (and only one) name and first name

- How can you make something like that unique?

- "When creating metadata for a resource, data categories are utilized for providing categories for metadata. For example the data category *resource title* is filled with the title of a resource and *metadata creator* contains the name of the creator of the metadata."

  http://www.sfs.uni-tuebingen.de/nalida/de/doku/glossar.html

- Attributes in AV-Structures

- Terminology from data base development

- Title
- Author
- Description
- Language
- Size
- Copyright
- Summary
- keyword(s)
- Address

- Institution
- Funder
- Subject field
- Type of resource
- Genre
- Programming language
- ...

- Select a resource you know or one that you created. Which data categories would you use to describe them?

- Which data categories would be inadquate for this type of resource?

- Could you group the data categories? Which data categories should always be grouped together?

CLARIN-D

- Wählen Sie eine Ressource, die Sie kennen oder erstellt haben. Welche Datenkategorien würden Sie verwenden, um Ihre Ressource zu beschreiben?
    – Siehe vorherige Folien

- Welche Datenkategorien würden für diese Ressource nicht zutreffen?
    – Für ein Korpus ist "Programmiersprache" nicht relevant, für eine "Software" vielleicht das Genre, etc.

- Können Sie Gruppen von Datenkategorien bilden? Welche Datenkategorien sollten immer zusammen verwendet werden?
    – Straße und Hausnummer, PLZ und Ort, Name und Rolle (Funktion), ....

# Chapter 6

Central registry of data categories

# Central registry of data categories

- Standardising of terminology
- Avoiding *Tag Abuse* by common definitions
  - Unique references by identifiers
  - Partially controlled vocabulary
  - Semantic interoperability: Referring to the semantics of a category
- Candidates for resources: Data categories created by others could be relevant for your own metadata
- Sample registry: CLARIN Concept registry
  - Partly replaces isocat.org ISO 12620:2009

www.isocat.org/interface/index.html

Google

Welcome Guest    Help

enter keywords here

Metadata

| # | Name | Version | Administration stat | Registration status | Check | Type | Scope |
|---|------|---------|--------------------|--------------------|-------|------|-------|
| | antecedent variable | 1:0 | private | private | ✓ | simple | public |
| 2623 | anthropological linguistics | 1:0 | private | private | ✓ | simple | public |
| 3788 | api | 1:0 | private | private | ✓ | open | public |
| 3786 | application type | 1:0 | private | private | ✓ | closed | public |
| 2624 | applied linguistics | 1:0 | private | private | ✓ | simple | public |
| 3787 | approach | 1:0 | private | private | ✓ | open | public |
| 4387 | archived | 1:0 | private | private | ✓ | simple | public |
| 4468 | artificial | 1:0 | private | private | ✓ | simple | public |
| 2653 | audio | 1:0 | private | private | ✓ | simple | public |
| 2689 | audio file format | 1:0 | private | private | ✓ | open | public |
| 4532 | audio recording | 1:0 | private | private | ✓ | simple | public |
| 4064 | audio transcription | 1:0 | private | private | ✓ | simple | public |
| 4115 | **author** | 1:0 | private | private | ✓ | open | public |
| 2453 | availability | 1:0 | private | private | ✓ | open | public |

My Workspace
- Public
  - Thematic Views
    - Metadata
      - Metadata
    - Morphosyntax
    - Semantic Content Repr
    - Syntax
    - Language Resource On
    - Lexicography
    - Language Codes
    - Terminology
    - Multilingual Informatio
    - Lexical Resources
    - Lexical Semantics
    - Translation
    - Sign language
    - Audio
  - Athens Core
  - BAS isocat group
  - CLARIN-NL/VL
  - Edisyn
  - GOLD
  - GilAndDan
  - GilAndSueEllen
  - KSU_PhD_Students
  - NKJP
  - RELISH
  - STTS
  - TBX-Basic
  - TBX-DCS

author - 1:0

## 2. Description Section

| Profile | Metadata |
|---------|----------|
| **2.1 Data Element Name Section** | |
| Data Element Name | author |
| Source | CMDI |
| [–]  **2.2 English Language Section** | |
| Language | English (en) |
| 2.2.1 Name Section | |
| Name | author |
| Name Status | preferred name |
| 2.2.2 Definition Section | |
| Definition | Indication of the author of, for instance, a book, a short story, a poem, or another piece of writing. |
| Source | NaLiDa |

- Go to the CLARIN Concept

    - https://openskos.meertens.knaw.nl/ccr/browser/

- Select the metadata section

- Find 10 data categories unfamiliar to you and find out what they mean

- What kind of resource could make use of such a description?

- Gehen Sie zum Datenkategorie Repository ISOcat

  – www.isocat.org

- Wählen Sie den Bereich Metadaten aus

- Finden Sie 10 Datenkategorien, die Ihnen nichts sagen und finden Sie heraus, was sich dahinter verbirgt

  – http://www.isocat.org/datcat/DC-2507 : Definition unter dem Link

  – http://www.isocat.org/datcat/DC-4507 …

# Data Category: author

| Key | 4115 |
|---|---|
| PID | http://www.isocat.org/datcat/DC-4115 |
| Type | complex/open |
| Owner | Hoppermann, Christina |
| Scope | public |

Unique idenfiter in ISOcat

## 1. Administration Information Section

### 1.1 Administration Record

| Identifier | author |
|---|---|
| Version | 1:0 |
| Registration Status | private |
| Administration Status | private |
| Justification | common metadata data category |
| *1.1.1 Creation* | |
| Creation Date | 2011-09-05 |
| Change Description | Definition of a new data category. |
| *1.1.2 Last Change* | |
| Last Change Date | 2012-03-19 |
| Change Description | Definition has been changed. |

## 2. Description Section

| Profile | Metadata |
|---|---|
| **2.1 Data Element Name Section** | |
| Data Element Name | author |
| Source | CMDI |
| **[−]  2.2 English Language Section** | |
| Language | English (en) |
| *2.2.1 Name Section* | |
| Name | author |
| Name Status | preferred name |
| *2.2.2 Definition Section* | |
| Definition | Indication of the author of a piece of writing. |
| Source | NaLiDa |

## 3. Conceptual Domain

## Data Category: author

| Key | 4115 |
|---|---|
| PID | http://www.isocat.org/datcat/DC-4115 |
| Type | complex/open |
| Owner | Hoppermann, Christina |
| Scope | public |

### 1. Admin

Unique name

#### 1.1 Administration Record

| Identifier | author |
|---|---|
| Version | 1:0 |
| Registration Status | private |
| Administration Status | private |
| Justification | common metadata data category |
| *1.1.1 Creation* | |
| Creation Date | 2011-09-05 |
| Change Description | Definition of a new data category. |
| *1.1.2 Last Change* | |
| Last Change Date | 2012-03-19 |
| Change Description | Definition has been changed. |

### 2. Description Section

| Profile | Metadata |
|---|---|
| **2.1 Data Element Name Section** | |
| Data Element Name | author |
| Source | CMDI |
| **[−]  2.2 English Language Section** | |
| Language | English (en) |
| *2.2.1 Name Section* | |
| Name | author |
| Name Status | preferred name |
| *2.2.2 Definition Section* | |
| Definition | Indication of the author of a piece of writing. |
| Source | NaLiDa |

### 3. Conceptual Domain

## Data Category: author

| Key | 4115 |
|---|---|
| PID | http://www.isocat.org/datcat/DC-4115 |
| Type | complex/open |
| Owner | Hoppermann, Christina |
| Scope | public |

### 1. Administration Information Section

| **1.1 Administration Record** | |
|---|---|
| Identifier | author |
| Version | 1:0 |
| Registration Status | private |
| Administration Status | private |
| Justification | common metadata data category |
| *1.1.1 Creation* | |
| Creation Date | 2011-09-05 |
| Change Description | Definition of a new data category. |
| *1.1.2 Last Change* | |
| Last Change Date | 2012-03-19 |
| Change Description | Definition has been changed. |

### 2. Description Section

| Profile | Metadata |
|---|---|
| **2.1 Data Element Name Section** | |
| Data Element Name | author |
| Source | CMDI |
| **[−]  2.2 English Language Section** | |
| Language | English (en) |
| *2.2.1 Name Section* | |
| Name | author |
| Name Status | preferred name |
| *2.2.2 Definition Section* | |
| Definition | Indication of the author of a piece of writing. |
| Source | NaLiDa |

English Definition

### 3. Conceptual Domain

## 2.2.2 Definition Section

| Definition | Indication of the author of a piece of writing. |
|---|---|

- Definitions in CCR on metadata define the category
  - Which bits of information can you find in the field of a data category?
  - → "Contains author"
  - Which value sets can you find there?
  - Where can you use this data category?
  - → Here: Example resource for which the data category is relevant

- Definitions **do not define** the semantics of a word
  - Here: they do not define what an author is

- Analytic definitions according to ISO 704
  - Consist of one fragment
  - Begin with the superordinate concept of the data category that is to be defined: direct superordinate term or above in a concept hierarchy
  - List essentially differentiating characteristics to related concepts
  - Do not contain a concept system
    - Concept-System and hierarchy are externally defined
    - Unification of definitions outside of the borders of a subject field
    - As general as possible

- Recycling of existing data categories if:
  - The concept in the CCR is the same as the needed one (the definition is sufficiently similar)
  - The definition is too narrow, but the creator or owner is willing to adjust the definition
  - The data category is from a different subject field but the meaning is obviously the desired one
  - The name of the data category is not the preferred variant

- Create a new Data Category if:
  - There is a different concepts/mental image/...
  - There is no similar term

- ## No, because
  - ### Persistency
  - ### Restricted creation
  - ### Only national content experts allowed

# Chapter 7

Generalizing groups of data categories

# Components

- Background
  - Use of similar data categories for different types of resources
  - Some data categories are typically used together
  - Selection of data categories is accelerated if they are grouped
- Groups of data categories: components
  - Cardinality: how often can or must a data category occur?
  - Value set: is there a closed vocabulary or a specific type of data?
- Components contain:
  - Data categories as elements
  - Other components
- Core model of the Component Metadata Infrastructure (CMDI, ISO 24622-1)
  - Co-developed by CLARIN
  - See: http://www.clarin.eu/cmdi

# Sample Components

- **General Info Component:**
  - Elements:
    - ResourceName
    - ResourceTitle
    - ResourceClass
    - Version
    - ...
  - Components:
    - Location
      - Address
      - Region
    - Descriptions
      - Description (in one language)

- **Access Component:**
  - Elements:
    - Availability
    - Distribution Medium
    - Catalogue Link
    - Price
  - Components:
    - DeploymentToolInfo
      - DeploymentTool
      - ToolType
      - ...
    - Contact
      - Person
      - Role
      - Address
      - ...

# Representing components: CMDI Component Specification Language (CCSL)

- Description language for components

- Defined in XML

- Descrives:
  - Elements of a component
    - cardinality
    - Reference to concept
    - Controlled vocabulary/data type
    - Multilinguality
    - Element name
  - &lt;CMD_Element Multilingual="true" CardinalityMax="1" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2453" name="Availability"/&gt;

# UML for CCSL

# Example CCSL in XML

CLARIN-D

```
1    <?xml version="1.0" encoding="UTF-8" standalone="yes"?>
2    <CMD_ComponentSpec isProfile="false" xsi:schemaLocation="http://www.clarin.eu/cmd http://www.clarin.eu/cmd/general-component-schema.xsd" xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
3        <Header>
4            <ID>clarin.eu:cr1:c_1290431694501</ID>
5            <Name>Access</Name>
6            <Description>Component contains information about the availability and accessibility of a resource.</Description>
7        </Header>
8        <CMD_Component CardinalityMax="1" CardinalityMin="1" name="Access">
9
10           <CMD_Element Multilingual="true"  CardinalityMax="1" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2453" name="Availability"/>
11           <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2967" name="DistributionMedium"/>
12           <CMD_Element CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="anyURI" ConceptLink="http://www.isocat.org/datcat/DC-2969" name="CatalogueLink"/>
13           <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2460" name="Price"/>
14           <CMD_Component CardinalityMax="unbounded" CardinalityMin="0" ComponentId="clarin.eu:cr1:c_1290431694500" name="DeploymentToolInfo">
15               <CMD_Element Multilingual="true" CardinalityMax="1" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2514" name="DeploymentTool"/>
16               <CMD_Element Multilingual="true"  CardinalityMax="1" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-3810" name="ToolType"/>
17               <CMD_Element Multilingual="true" CardinalityMax="1" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2547" name="Version"/>
18               <CMD_Element CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="anyURI" ConceptLink="http://www.isocat.org/datcat/DC-63" name="Url"/>
19               <CMD_Component CardinalityMax="1" CardinalityMin="0" ComponentId="clarin.eu:cr1:c_1290431694486" name="Descriptions">
20                   <CMD_Element Multilingual="true" DisplayPriority="1" CardinalityMax="unbounded" CardinalityMin="1" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2520" name="Description">
21                       <AttributeList>
22                           <Attribute>
23                               <Name>type</Name>
24                               <ValueScheme>
25                                   <enumeration>
26                                       <item ConceptLink="">short</item>
27                                       <item ConceptLink="">long</item>
28                                   </enumeration>
29                               </ValueScheme>
30                           </Attribute>
31                       </AttributeList>
32                   </CMD_Element>
33               </CMD_Component>
34           </CMD_Component>
35           <CMD_Component CardinalityMax="unbounded" CardinalityMin="1" ComponentId="clarin.eu:cr1:c_1290431694487" name="Contact">
36               <CMD_Element Multilingual="false" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2978" name="Person"/>
37               <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-3807" name="Role"/>
38               <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2505" name="Address"/>
39               <CMD_Element Multilingual="false" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2521" name="Email"/>
40               <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-3812" name="Department"/>
41               <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2979" name="Organisation"/>
42               <CMD_Element Multilingual="false" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2461" name="TelephoneNumber"/>
43               <CMD_Element Multilingual="false" CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2455" name="FaxNumber"/>
44               <CMD_Element CardinalityMax="unbounded" CardinalityMin="0" ValueScheme="anyURI" ConceptLink="http://www.isocat.org/datcat/DC-63" name="Url">
45                   <AttributeList>
```

Bundesministerium
für Bildung
und Forschung

- What is a component

- What kind of information can you find in a component?

- How does the Concept Registry relate to components?

- Advanced issues:

  - What do you need to create components?

**CLARIN-D**

- Was ist eine Komponente?
  - Eine Gruppierung von Datenkategorien
- Welche Informationen findet man in einer Komponente?

  - Welche Elemente und Komponenten Teil der Komponente sind
  - Häufigkeit
  - ISOcat-Referenz
  - Kontrolliertes Vokabular/Datentyp
  - Mehrsprachigkeit
  - Elementname
- Was hat ISOcat mit Komponenten zu tun?
  - Die Elementbeschreibungen kommen aus ISOcat
- Weiterführende Fragen:
  - Was brauchen Sie, um selbst eine Komponente zu erstellen?
    - Werkzeuge, die die Komplexität entfernen

# Kapitel 8

Generalising for types of resources

- Background
  - Each type of resource receives a metadata schema
  - Reuse of components
  - Different types of resources require different schemas
- One type of resource = one profile
  - Basis for generating an XSchema for validation
- Profile description contains
  - Data categories as elements
  - Additional components
- Core model of the Component Metadata Infrastructure (CMDI)
  - Developed together with the CLARIN
  - See: http://www.clarin.eu/cmdi

- Dublin Core:
  - Dublin Core Metadata Initiative (DCMI)
  - Cataloguing of documents
  - Originally standardised DC: 15 core elements
  - qualified DC: further specification
- OLAC: Open Language Archives Community
  - Aiming at a library of language resources
  - Extending Dublin Core (DC)
  - Contains DC
  - More language related stuff
- IMDI: ISLE MetaData Initiative
  - Metadata for language resources
  - Especially spoken language
  - Very rich descriptions
  - Including speaker, roles, age, ….

- Experiments
  - Soziolinguistics
  - Psycholinguistics
- Lexical resourcen
- Diachrone textccorpora
- Corpora of spoken language
- Corpora of written language
- Software Tools
- Webservices
- Webservice-Toolchain

- Profiles defined as components
- Specified in CCSL: component type profile

```xml
<?xml version="1.0" encoding="UTF-8" standalone="yes"?>
<CMD_ComponentSpec isProfile="true"
    xsi:schemaLocation="http://www.clarin.eu/cmd http://www.clarin.eu/cmd/general-component-schema.xsd"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance">
    <Header>
        <ID>clarin.eu:cr1:p_1290431694580</ID>
        <Name>TextCorpusProfile</Name>
        <Description>A CMDI profile for text (i.e. written) corpus resources.</Description>
    </Header>
    <CMD_Component CardinalityMax="1" CardinalityMin="1" name="TextCorpusProfile">
        <CMD_Component CardinalityMax="1" CardinalityMin="1"
            ComponentId="clarin.eu:cr1:c_1290431694495" name="GeneralInfo">
            <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0"
                ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2544"
                name="ResourceName"/>
            <CMD_Element Multilingual="true" CardinalityMax="unbounded" CardinalityMin="0"
                ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2545"
                name="ResourceTitle"/>
            <CMD_Element DisplayPriority="1" CardinalityMax="unbounded" CardinalityMin="1"
                ConceptLink="http://www.isocat.org/datcat/DC-3806" name="ResourceClass">
                <ValueScheme>
                    <enumeration>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4360">Lexicon</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4361">Corpus</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4362">Tool</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4363">Grammar</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4364">Fieldwork Material</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4365">Experimental Data</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4366">Survey Data</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4367">Test Data</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4368">Toolchain</item>
                        <item AppInfo="" ConceptLink=" http://www.isocat.org/datcat/DC-4369">ResourceBundle</item>
                        <item AppInfo="" ConceptLink="http://www.isocat.org/datcat/DC-2599">other</item>
                        <item AppInfo="" ConceptLink="http://www.isocat.org/datcat/DC-2591">unknown</item>
                    </enumeration>
                </ValueScheme>
            </CMD_Element>
            <CMD_Element Multilingual="true" CardinalityMax="1" CardinalityMin="0"
                ValueScheme="string" ConceptLink="http://www.isocat.org/datcat/DC-2547"
                name="Version"/>
            <CMD_Element CardinalityMax="1" CardinalityMin="0"
                ConceptLink="http://www.isocat.org/datcat/DC-3818" name="LifeCycleStatus">
                <ValueScheme>
                    <enumeration>
```

- What is a profile?

- What is in a profile?

- What is the connection between profiles and components?

- Advanced question:
  - What do you need to create a profile?

- ## Was ist ein Profil?
  - Eine besondere Komponente für einen bestimmten Typ von Ressource

- ## Welche Informationen findet man in einem Profil?
  - Alles, was man auch in einer Komponente findet plus die Information, dass es sich um ein Profil handelt

- ## Was ist der Zusammenhang von Profil und Komponente?
  - Profile sind Komponenten und enthalten in der Regel Komponenten und Elemente

- ## Weiterführende Fragen:
  - Was brauchen Sie, um selbst ein Profil zu erstellen?
    - Immer noch Werkzeuge mit geringerer Komplexität

# Chapter 9

## CMDI Component Registry

- Die Component Registry ist eine Webapplikation, mit der man sich vor dem Kurs vertraut machen sollte.

- Zugang: http://catalog.clarin.eu/ds/ComponentRegistry/#

- Zum Erstellen von Profilen und Komponenten ist ein CLARIN-Webseiten-Account erforderlich. Wenn Sie nicht Partner von CLARIN sind, beachten Sie bitte die Bedingungen, unter denen Sie Zugang erhalten, siehe http://www.clarin.eu/user

- Die Erfahrung zeigt, dass schwer vorherzusagen ist, wie lange Nutzer brauchen, um sich mit neuen Applikationen vertraut zu machen. Möglicherweise ist es notwendig, dass Sie diese Schritte vorführen.

- High complexity of component model

- Tools not part of the models

  – Creation of components

  – Connecting components with components

  – Include data categories

- Reuse of components requires central registry

  – Unique reference

  – Easily accessible

**Clarin Component Browser**

Browse...   Edit...   Import...

**Profiles** | **Components**

Public space ▾ | 🔍 filter... | Showing 52 of 52

| Name | Group Nam | Domain Name | Creator | Description | Registration Date |
|------|-----------|-------------|---------|-------------|-------------------|
| AnnotationTool | | | Eric Sa... | Description of a tooladapted from adept for interviews | 09 February 2011 14:13:59 |
| BamdesLexicalResource | | computational_linguistics | Dieter V... | Lexical Resource as used by BAMDES (for theharvestingday.eu) | 27 October 2010 15:47:42 |
| BamdesMultimodalCorpus | | computational_linguistics | Dieter V... | Oral Corpus as used by BAMDES (for theharvestingday.eu) | 27 October 2010 16:00:41 |
| BamdesOralCorpus | | computational_linguistics | Dieter V... | Oral Corpus as used by BAMDES (for theharvestingday.eu) | 27 October 2010 16:00:05 |
| BamdesTool | | computational_linguistics | Dieter V... | Tool as used by BAMDES (for theharvestingday.eu) | 27 October 2010 15:48:54 |
| BamdesWrittenCorpus | | computational_linguistics | Dieter V... | Written Corpus as used by BAMDES (for theharvestingday.eu) | 27 October 2010 15:59:28 |
| Bedevaartbank | | | Folkert ... | Profile for Bedevaartbank | 30 July 2010 9:33:53 |
| Boedelbank | | | Folkert ... | Profile for Boedelbank | 30 July 2010 9:36:25 |
| cmdi-virtual-collection | | | Patrick ... | A component for describing personalized sets of metadata descriptions | 21 April 2010 16:18:30 |
| collection | | | Matej D... | minimal collection profile (started for OLAC records) | 15 October 2010 20:58:37 |
| component-dc-terms | | | Patrick ... | DC metadata | 21 April 2010 16:19:45 |
| component-dc-terms-modular | | | Patrick ... | DC metadata | 21 April 2010 16:19:45 |
| DBNL | | | Folkert ... | Profile for DBNL | 30 July 2010 9:37:50 |
| DcmiTerms | | Other | Dieter V... | DCMI Terms vocabulary, see http://dublincore.org/documents/dcmi-terms/ | 28 October 2010 15:41:46 |
| DIDDD | | | Folkert ... | Profile for DIDDD | 30 July 2010 9:39:00 |

# Clarin Component Browser: Overview of existing public components and profiles

- Profile name

- Description

- Date

- Creator

- Group

- Profile-URL→ Show Info

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

# Creating new profiles and components

- Elements referring to Concept Registry
- Selecting compone...
- Cardinality
- Data types for elem...

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

- Go to http://catalog.clarin.eu/ds/ComponentRegistry/
  - Click "login"
  - Login using Shibboleth
- Explore a couple of profiles
- Click edit
- Create some profile and components; do NOT publish them
- Tutorial http://www.sfs.uni-tuebingen.de/nalida/en/docu/tutorials/component-registry.html

- ## Advantage reuse:
  - Quickly create components and profiles
  - Components are "better" testet and probably more complete
  - Reuse of tools
    - XPaths with a couple of wildcards
    - No need for individual applications

- ## Advantage new creation
  - More of a control

- ## Recommendation: Reuse whenever possible or create new ones on the basis of existing ones

# Chapter 10

## Creating CMDI instances

- Three sections
  - Header: Description of the metadata
  - Resources: reference to the resources

- Components: according to the type of resource

```
          CMD
         / | \
   Header Resources Components
```

# CMDI Header

- Information:
  - Name of metadata creator: MdCreator
  - Date of metadata creation: MdCreationDate
  - URL of metadata file: MdSelfLink
  - Identifier of the CMDI Profile: MdProfile
  - Name of the archive: MdCollectionDisplayName
- Creation:
  - Archiving process
  - URL of the metadaten file persistent
    - Name of the creator as known by the (Login)
    - CMDI profile known by selection
    - Name of Archive known by archive
    - Creation without archiving process: manual or by defaults

- Information
  - ResourceProxyList: List of PIDs of primary data JournalFileProxyList: List of files with provenance metadata
  - ResourceRelationList: Relation to other resources
  - IsPartOfList: if the resource is part of a larger unit
- Creation:
  - Archiving process:
    - ResourceProxyList: PIDs created during archiving
  - Other lists as needed

- ## Information:
  - Components
  - Data categories from Concept Registry
- ## Creation:
  - In  an editor
    - XML editor
    - Text editor
    - Metadata editor
  - By transformation
    - From metadata in other formats
    - From relational databases

```xml
127    <Components>
128        <ExperimentProfile>
129            <GeneralInfo ComponentId="clarin.eu:cr1:c_1290431694495">
130                <ResourceName>Stroop Nouns Bewegung</ResourceName>
131                <ResourceTitle>Stroop Nouns Bewegung</ResourceTitle>
132                <ResourceClass>Experimental Data</ResourceClass>
133                <Version/>
134                <LifeCycleStatus>archived</LifeCycleStatus>
135                <StartYear>2010</StartYear>
136                <CompletionYear>2010</CompletionYear>
137                <LegalOwner>Martin Lachmair</LegalOwner>
138                <Location ComponentId="clarin.eu:cr1:c_1290431694494">
139                    <Address>Friedrichstr. 21, 72070 Tübingen</Address>
140                    <Country ComponentId="clarin.eu:cr1:c_1290431694493">
141                        <CountryName xml:lang="en">Germany</CountryName>
142                        <CountryCoding>DE</CountryCoding>
143                    </Country>
144                </Location>
145                <Descriptions ComponentId="clarin.eu:cr1:c_1290431694486">
146                    <Description xml:lang="en">This resource provides a reaction time experiment within the context of the
147                        SFB 833 project B4.</Description>
148                </Descriptions>
149            </GeneralInfo>
150            <Project ComponentId="clarin.eu:cr1:c_1290431694522">
151                <ProjectName>B4</ProjectName>
152                <ProjectTitle xml:lang="de">Simulationsansatz des Sprachverstehens: Komposition von Bedeutung</ProjectTitle>
153                <ProjectTitle xml:lang="en">The experiential simulation view of language comprehension: How is sentence meaning composed?</ProjectTitle>
154                <ProjectID/>
155                <Funder xml:lang="de">Deutsche Forschungsgemeinschaft (DFG)</Funder>
156                <Funder xml:lang="en">German Research Foundation (DFG)</Funder>
157                <Url>http://www.sfb833.uni-tuebingen.de/b-bereich-kognition/b4-kaup.html</Url>
158                <Institution ComponentId="clarin.eu:cr1:c_1290431694496">
159                    <Department xml:lang="de">Sonderforschungsbereich 833: Bedeutungskonstitution-
160                        Dynamik und Adaptität sprachlicher Strukturen</Department>
161                    <Department xml:lang="en">Collaborative Research Centre 833: The construction of meaning - the dynamics
162                        and adaptivity of linguistic structures</Department>
163                    <Organisation xml:lang="de">Eberhard Karls Universität Tübingen</Organisation>
164                    <Organisation xml:lang="en">University of Tübingen</Organisation>
165                    <Url>http://www.sfb833.uni-tuebingen.de/</Url>
166                    <Location ComponentId="clarin.eu:cr1:c_1290431694494">
167                        <Address>Nauklerstraße 35, 72074 Tübingen</Address>
168                        <Region>Baden-Württemberg</Region>
169                        <Country ComponentId="clarin.eu:cr1:c_1290431694493">
170                            <CountryName xml:lang="de">Deutschland</CountryName>
171                            <CountryName xml:lang="en">Germany</CountryName>
172                            <CountryCoding>DE</CountryCoding>
173                        </Country>
```

XPath 2.0

Text  Raster  Autor

# Creation of a CMDI instance in an XML-Editor

- Edit in XML code
  - Syntax highlighting
  - automatic validation

- disadvantage:
  - complexity
  - Learning curve of archivist

# Comedi

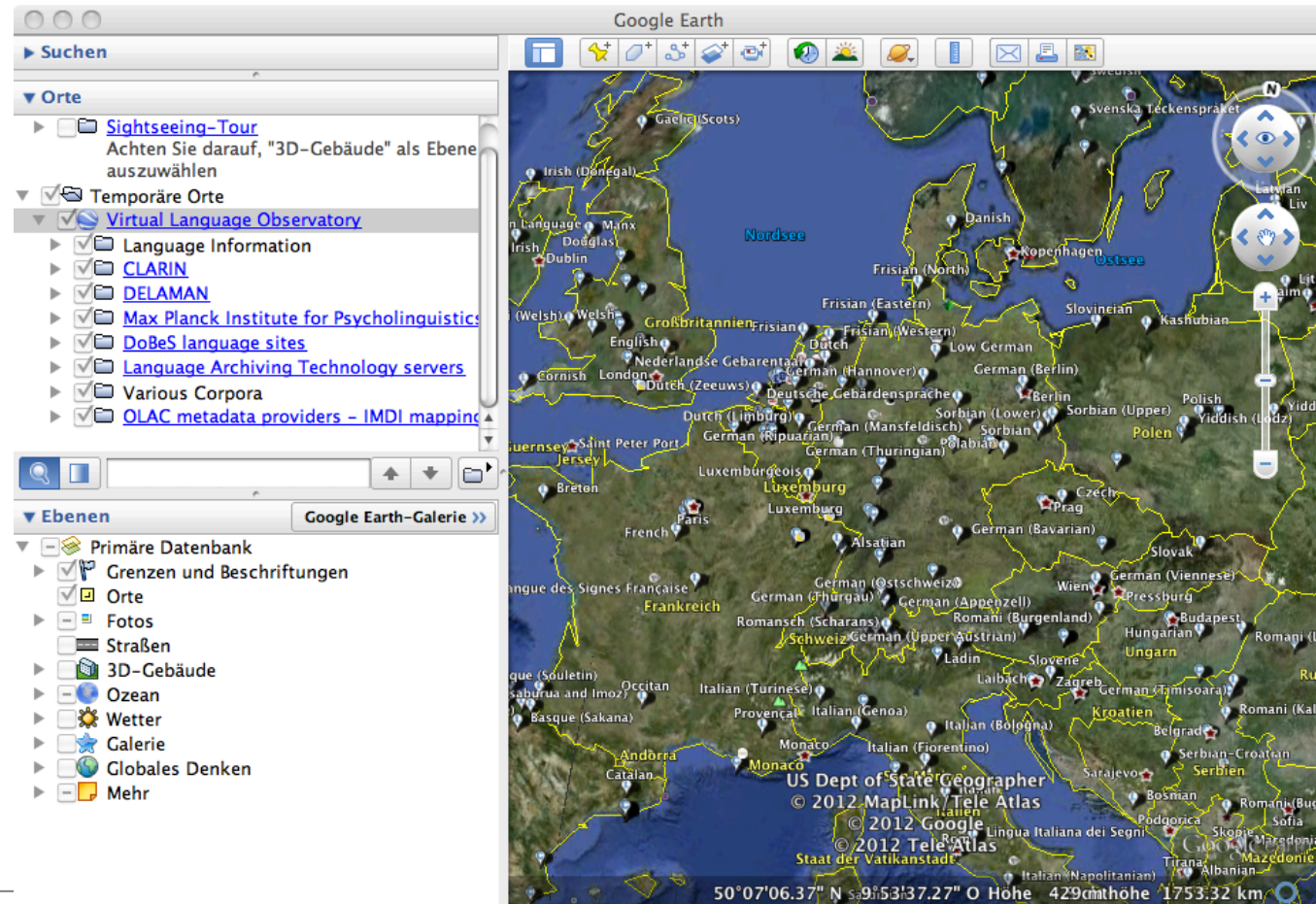- http://clarino.uib.no/comedi/metadata-editor

# Chapter 11

Using metadata

# Exchanging metadata

- Spreading the word for visibility
- Application/search engine for resources
- Resource as such can stay at repository
- Integration into other data sources
  - Library catalogues
  - Websites
- Protocol: OAI-PMH
  - Open Archives Initiative – Protocol for Metadata Harvesting
  - Metadata on OAI-PMH Server
  - Clients harvest changed data
- Most important:
  - Search

# Visualising resources in Google-Earth

- See http://www.clarin.eu/vlo/

# Virtual Language Observatory

**VLO**

Explore the world of language resources and technology from different perspectives

search

## COLLECTION

Tübingen Language Resources (19062)
childes (14429)
Language and Cognition (9910)
Endangered Languages (9600)
Ethnologue: Languages of the World (7413)
WALS RefDB (7348)
talkbank (6873)
MPI CGN (6725)
Acquisition (6009)
Oxford Text Archive (4598)
more...

## LANGUAGE

English (40705)
German (33664)
Dutch (14518)
Spanish; Castilian (6091)
French (4351)
Turkish (3289)
Japanese (3281)
Chinese (1189)
Polish (1149)
Indonesian (935)
more...

Showing 1 to 10 of 123524

<< < *1* 2 3 4 5 6 7 8 9 10 > >>

### Results

A letter from the Hon. Thomas Hervey to Sir Thomas Hanmer [Electronic resource]

ACL Membership List

Computational Linguistics, Volume 22, Number 4, December 1996

Tarlton's Newes out of Purgatorie [Electronic resource]

The tempest [Electronic resource], or, The enchanted island : a comedy, as it is now acted at his Highness the Duke of York's theatre

!O!ung: a language of Angola

!Xóõ: a language of Botswana

"Expertness" from Structured Text? RECONSIDER: A Diagnostic Prompting Program

"Le Monde Diplomatique" Arabic tagged corpus

"Le Monde Diplomatique" Text corpus in Arabic

## CONTINENT

Europe (32260)
North-America (11098)
Asia (8198)
South-America (3887)
Oceania (2377)
Africa (1536)
Middle-America (1197)
Australia (829)
North America (185)
Australien (2)
more...

## GENRE

discourse (37293)
referenz (3797)
legende (3598)
primary_text (3308)
gedicht (3294)
spontanous speech (3240)
language_description (2903)
stimuli (1577)
erzählung (1387)
fiktion (1323)
more...

## COUNTRY

Germany (10514)
United States (10088)
Netherlands (9675)
Japan (3269)
United Kingdom (2852)
Papua New Guinea (2732)

## SUBJECT

general_linguistics (5902)
typology (5896)
text_and_corpus_linguistics (5166)
syntax (4514)
semantics (2558)
monologue about free topic (2265)

# CLARIN Virtual Language Observatory

- Search by data categories
- Based on CMDI metadata

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

# Chapter 12

Summary

GEFÖRDERT VOM

Bundesministerium
für Bildung
und Forschung

- Description of resources

- Purpose: Overview and locating resources

- Different metadata formats

- Different resource types = different metadata schemas

- Integration of metadata in archive systems/repositories/catalogues

- Description by filling in data categories

- Definition of data categories: Concept Registry

- Grouping of data categories: components

- Components combined to describe one type of resource: profile

- Tools:

  - Concept Registry: Data categories

  - Component Registry: creation and central storage of components

  - ARBIL and Comedi: CMDI editors

- What is the connection between metadata, data categories, components and profiles?
- Why is it necessary to create metadata?
- Create metadata for a book using CMDI!
  - Select data categories
  - Create a component
  - Create an instance matching the component

Contact: Seminar für Sprachwissenschaft

**Dr. Thorsten Trippel**

Wilhelmstr. 19, 72074 Tübingen

Telefon: +49 7071 29-77352

thorsten.trippel@uni-tuebingen.de