



# Prozesse zur Beschreibung und Archivierung linguistischer Forschungsdaten

Christina Hoppermann, Thorsten Trippel, Claus Zinn  
Initiative zur Nachhaltigkeit Linguistischer Daten (NaLiDa)

Die Archivierung und Beschreibung von Forschungsdaten sind Aspekte, die nicht nur zur guten wissenschaftlichen Praxis zählen, sondern insbesondere von Forschungsförderungsorganisationen zunehmend über einen festgelegten Zeitraum, wie die von der Deutschen Forschungsgemeinschaft (DFG) geforderten 10 Jahre, vorgeschrieben werden. Diese Vorgaben stellen jedoch Herausforderungen an die Datenersteller, die über ihre fachwissenschaftliche Arbeit hinausgehen. Um diesen Mehraufwand minimieren zu können, stellt die Initiative zur Nachhaltigkeit Linguistischer Daten (NaLiDa) sowohl Verfahrensweisen als auch Anleitungen und Referenzen zur Beschreibung und Archivierung linguistischer Forschungsdaten zur Verfügung.

## Archivierung

### Wozu dienen Repositorien?

Ein Repository (Archiv) bezeichnet einen (digitalen) Ort zur nachhaltigen, langfristigen Sicherung von Daten und Dokumenten unter Einhaltung von Zugangsberechtigungen.

Repositorien bieten folgende Vorteile:

- Fähigkeit zur persistenten Datenspeicherung → Zugänglichkeit von Ressourcen,
- Zitierbarkeit von Ressourcen durch die Vergabe persistenter Identifikatoren (PIDs),
- Verbreitung von Ressourcen durch ihre Metadaten mittels OAI-PMH (Open Archives Initiative Protocol for Metadata Harvesting),
- Auffindbarkeit von Ressourcen durch Suchmöglichkeiten (z.B. durch einen Faceted Browser, der auf Metadaten basiert),
- Erfüllung der Vorgaben von Forschungsförderungsorganisationen bezüglich der Archivierung.

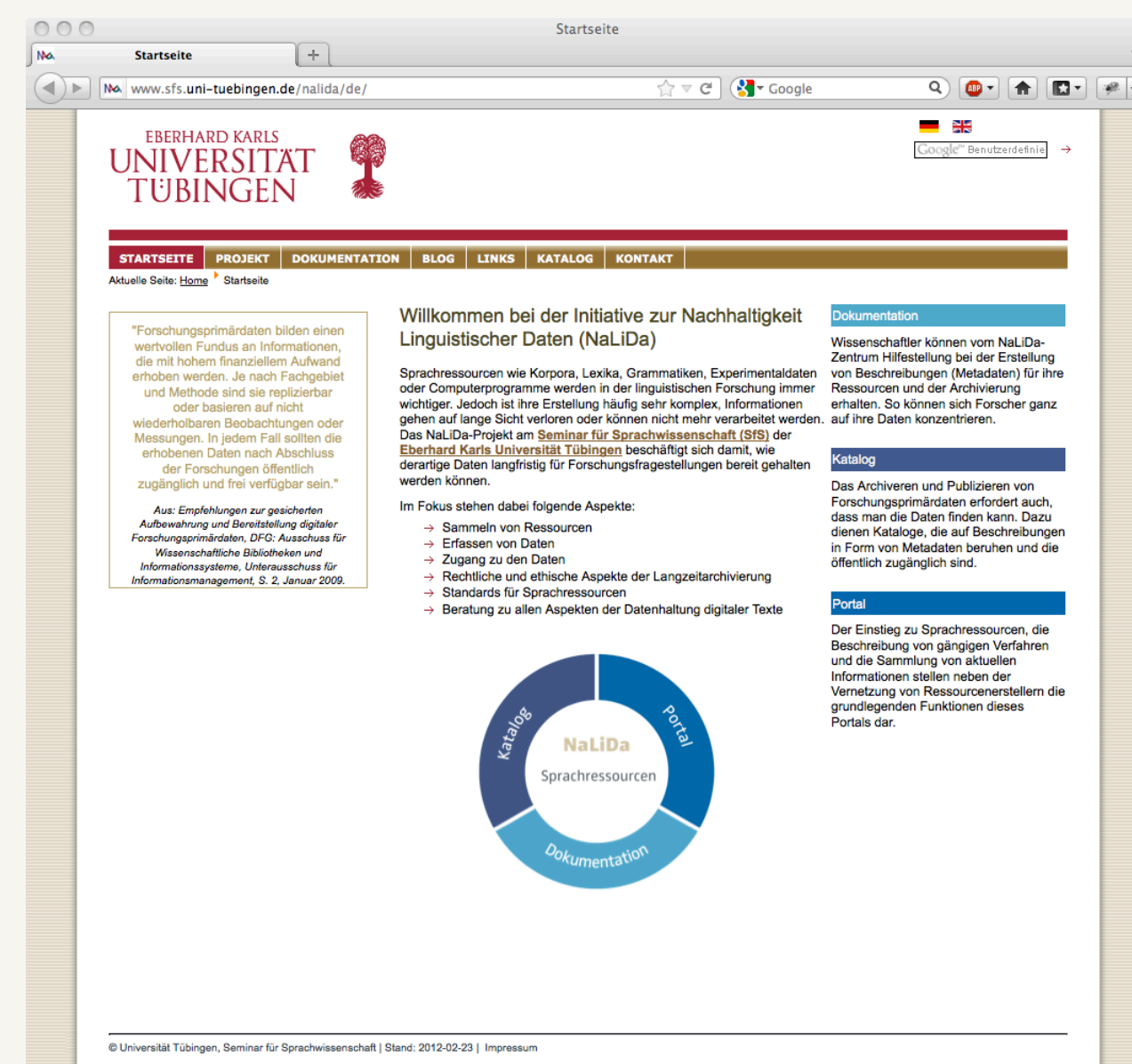
### Organisation in digitale Objekte

Repositoriensysteme wie Fedora-Commons ([www.fedora-commons.org](http://www.fedora-commons.org)) verwenden als zentrale Einheit *digitale Objekte*, die wiederum Datenströme enthalten.

- Digitales Objekt
  - Ressource mit ihren Metadaten und zugehörigen Dateien
  - Vergleichbar mit einem Ordner bzw. Verzeichnis bei Dateisystemen
  - Adressierbar mittels eines persistenten Identifikators (PID)
- Datenstrom:
  - Einzelner Bestandteil des digitalen Objektes
  - Vergleichbar mit Datei im Verzeichnis
- Vorteile dieses Ansatzes:
  - Nicht auf eine Ressourcenklasse beschränkt
  - Digitale Objekte beinhalten verschiedene Datentypen
  - Metadaten sind Datenstrom in dem Objekt
- Nachteile dieses Ansatzes:
  - Flache Struktur
  - Hierarchien nur durch Angabe von Relationen von Objekten
- Konsequenz: verschiedene Hierarchien sind gleichzeitig möglich



## NaLiDa bei der Beschreibung und Archivierung von Forschungsdaten



Sprachressourcenportal der Initiative zur Nachhaltigkeit Linguistischer Daten (NaLiDa): [www.sfs.uni-tuebingen.de/nalida](http://www.sfs.uni-tuebingen.de/nalida)

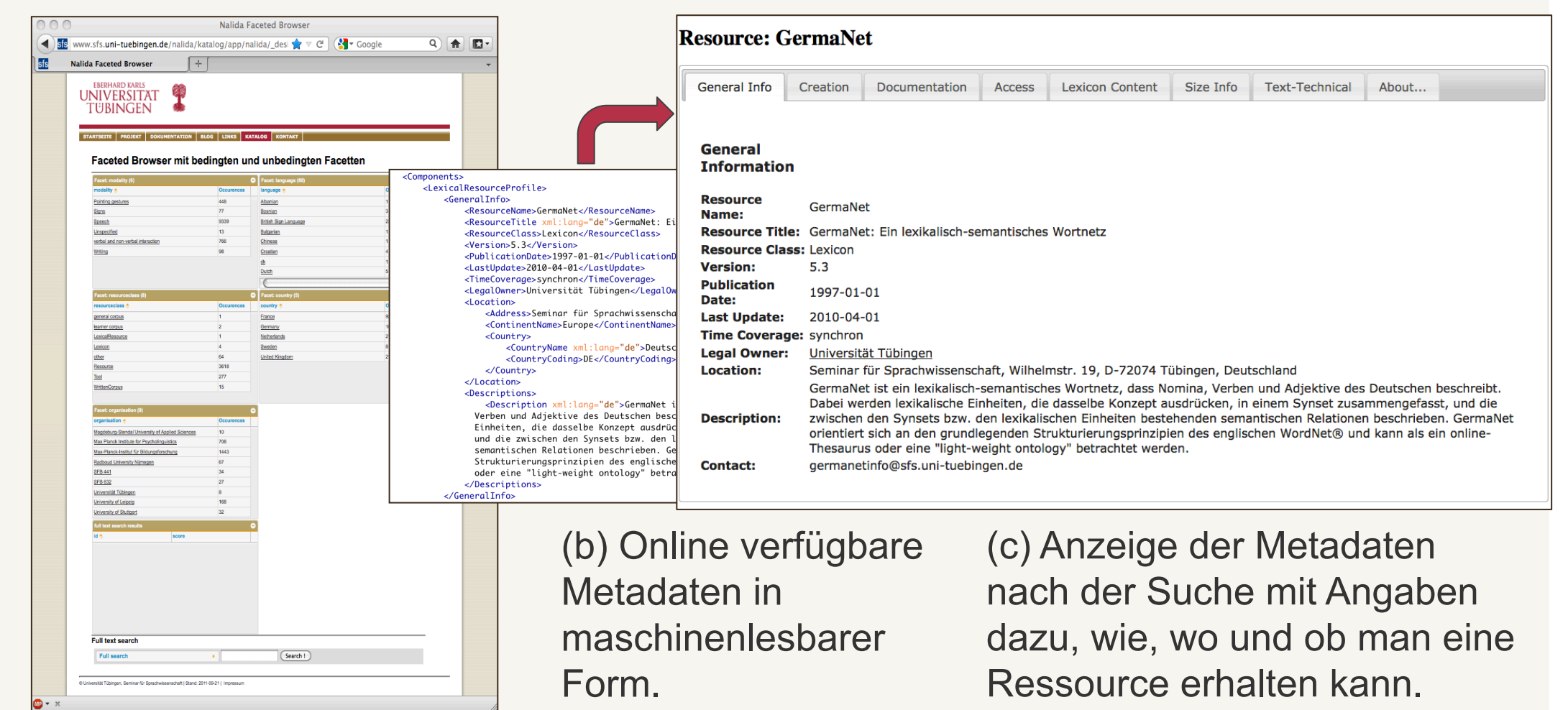
NaLiDa hat folgende Schwerpunkte:

- Unterstützung der Archivierung durch die Erstellung von Metadaten:
  - Bereitstellung von Profilen (Beschreibungsvorlagen) für verschiedene Arten von Ressourcen
  - Unterstützung bei der Verwendung der Profile
  - Bearbeitung existierender Bestandsmetadaten in anderen Formaten
  - Hilfe bei der Erstellung von Metadaten, z.B. durch den Metadaten-Editor ProFormA
- Bereitstellung von Informationen zum Archivieren von Forschungsprimärdaten:
  - Erstellung von digitalen Objekten (Identifikation der Bestandteile)
  - Beschreibung der Arbeitsabläufe zur Archivierung
  - Beratung zu nachhaltigen Datenformaten für Primärdaten
- Erfassen von Daten:
  - Unterstützung bei der Erfassung von Primärdaten und Metadaten
  - Suche über die Metadaten
- Standards für Sprachressourcen (insbesondere für Metadaten)

## Beschreibung der Daten: Metadatenerstellung

### Wozu dienen Metadaten?

- Metadaten sind strukturierte Daten, die sowohl zur Beschreibung als auch zum Auffinden von Forschungsdaten verwendet werden können
- Auffindbarkeit von Ressourcen
- Überblick über Ressourcen und Zugang zu Ressourcen
- Zitierbarkeit sowohl der Ressourcen selbst als auch der Metadaten

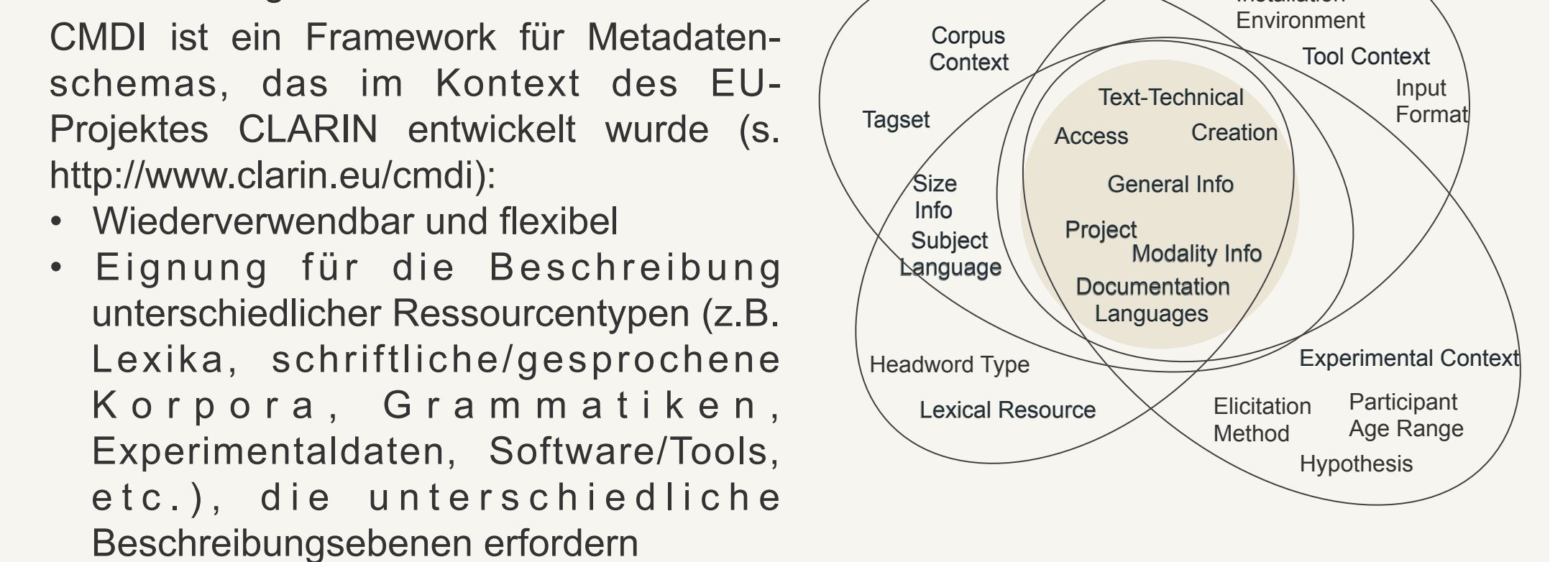


(a) Faceted Browser in Kombination mit Volltextsuche als Suchfunktion über Metadaten.

### Component MetaData Infrastructure (CMDI)

CMDI basiert auf drei grundlegenden Konzepten:

- Profile
  - Beschreibungsvorlagen für Ressourcen-typen
  - beinhalten Komponenten als Bausteine
- Komponenten
  - gruppieren Datenkategorien zu semantischen Einheiten
  - können selbst weitere Komponenten enthalten
  - wiederverwendbar in unterschiedlichen Profilen
- Datenkategorien

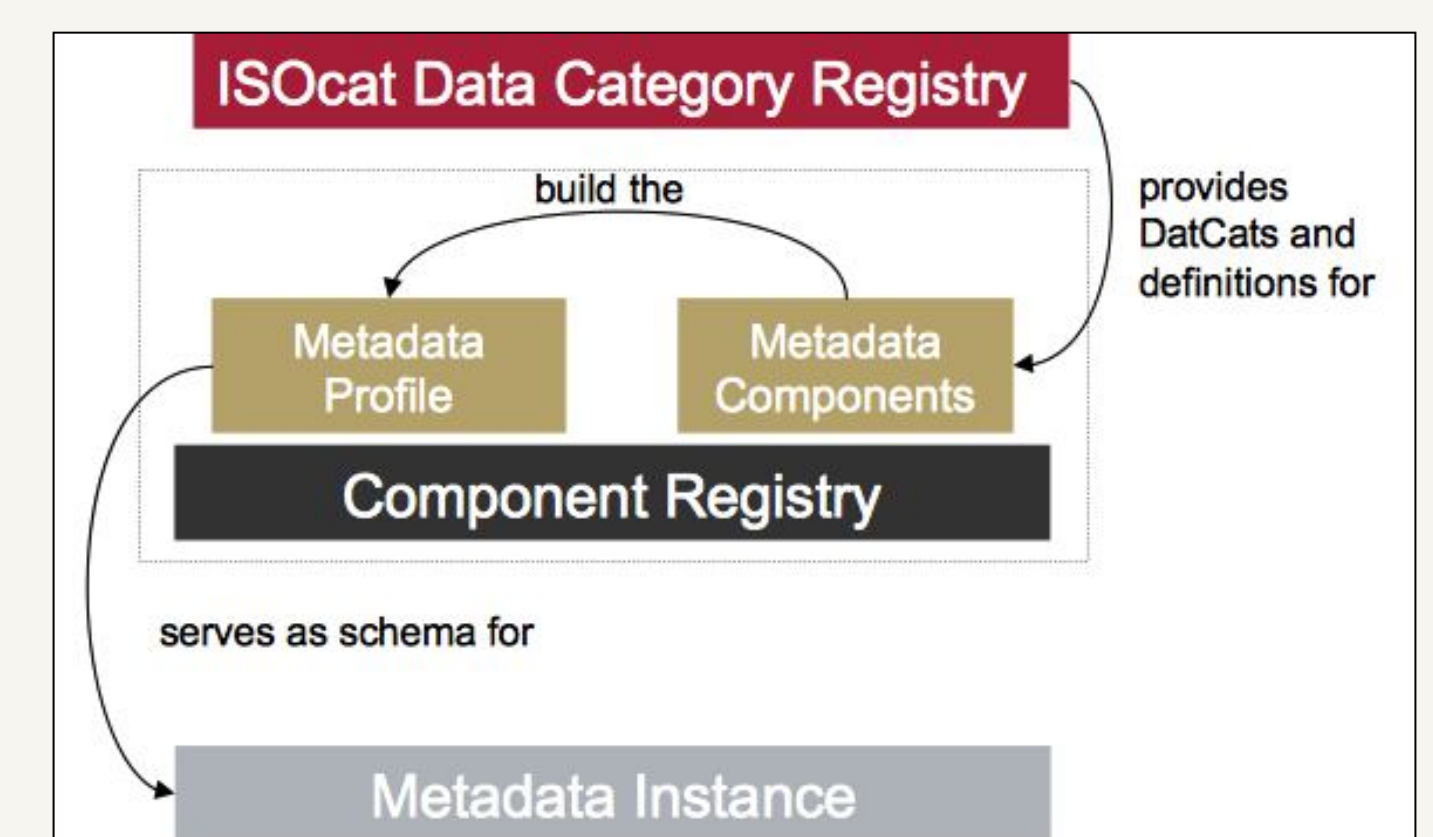


CMDI ist ein Framework für Metadaten-schemas, das im Kontext des EU-Projektes CLARIN entwickelt wurde (s. <http://www.clarin.eu/cmd/>):

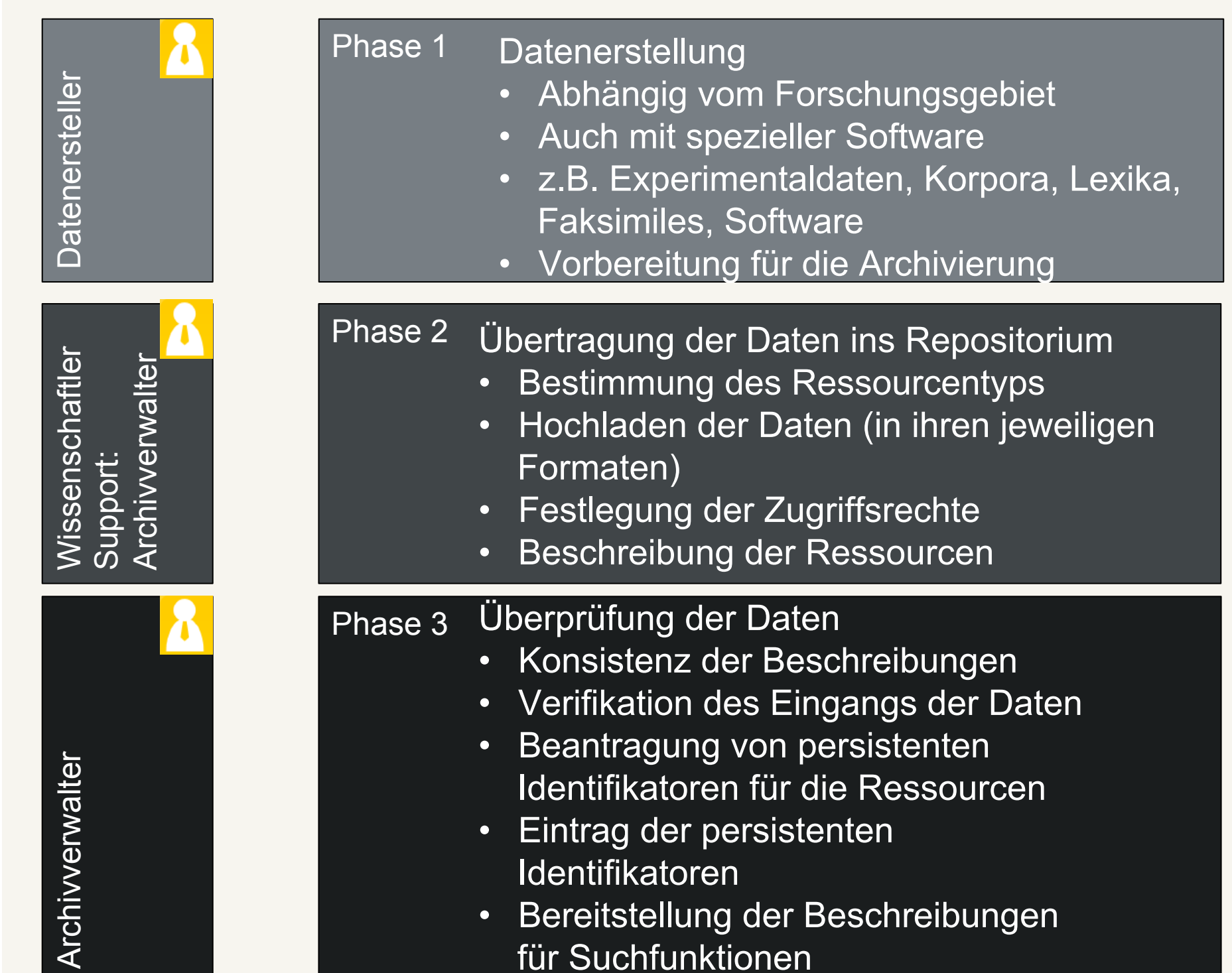
- Wiederverwendbar und flexibel
- Eignung für die Beschreibung unterschiedlicher Ressourcentypen (z.B. Lexika, schriftliche/gesprochene Korpora, Grammatiken, Experimentaldaten, Software/Tools, etc.), die unterschiedliche Beschreibungsebenen erfordern

Zentrale Verzeichnisse zur Definition von CMDI Beschreibungsmitteln:

- ISOcat (<http://www.isocat.org>)
- Component Registry (<http://catalog.clarin.eu/ds/ComponentRegistry/#>)



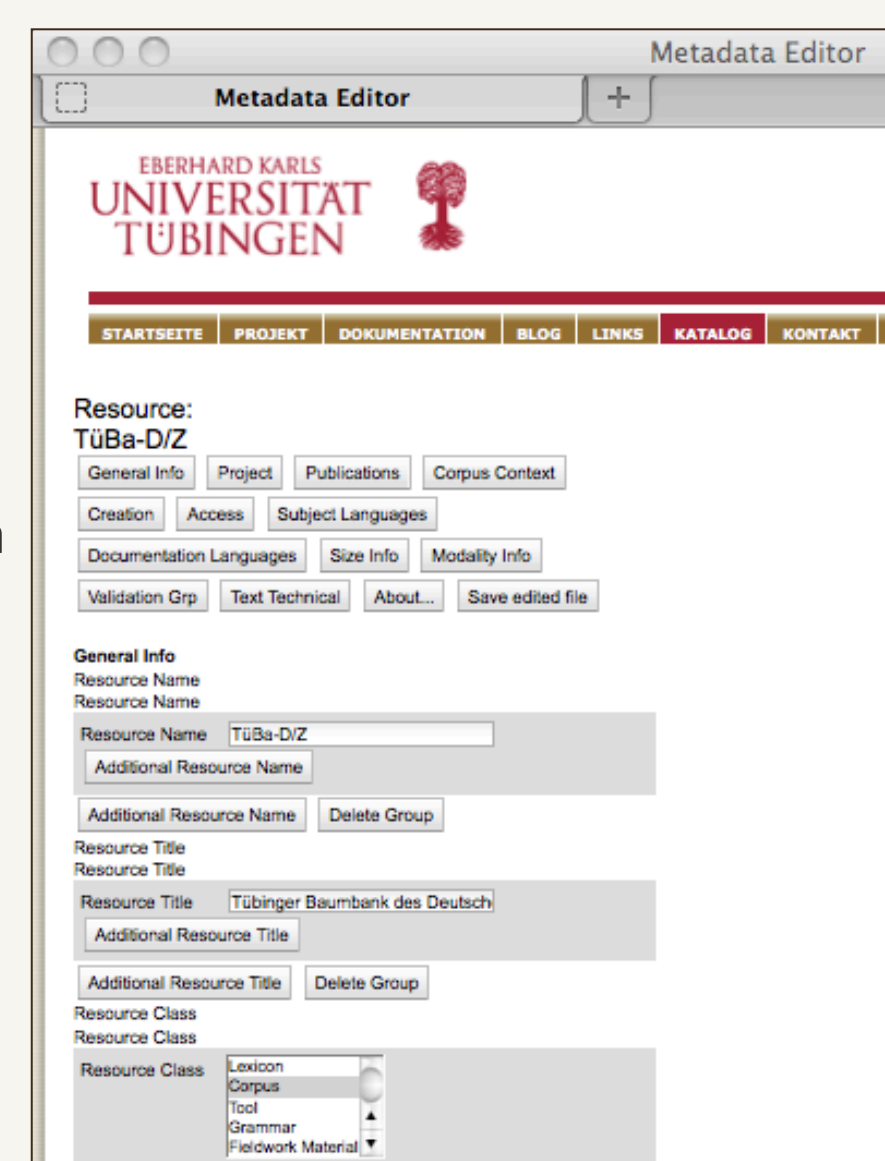
## 3 Phasen der Datenarchivierung



## ProFormA: Ein Editor zur Unterstützung bei der Erstellung von Metadaten

ProFormA:

- Formular- und webbasierter Editor
- Verwendung: Erstellen und Editieren von Metadaten
- Basiert auf XForms, einem W3C-Standard zur Erstellung webbasierter Formulare
- Repräsentiert zugrunde liegende XML-Strukturen von Metadaten in einer nutzerfreundlichen Formularansicht in HTML



Vorteile:

- Minimiert die Komplexität des Metadatenerstellungsprozesses
- Kann ohne (XML-)Vorkenntnisse angewendet werden
- Stellt keinen Mehraufwand für Nutzer dar

## Zusammenhang von Metadaten und Ressourcen

