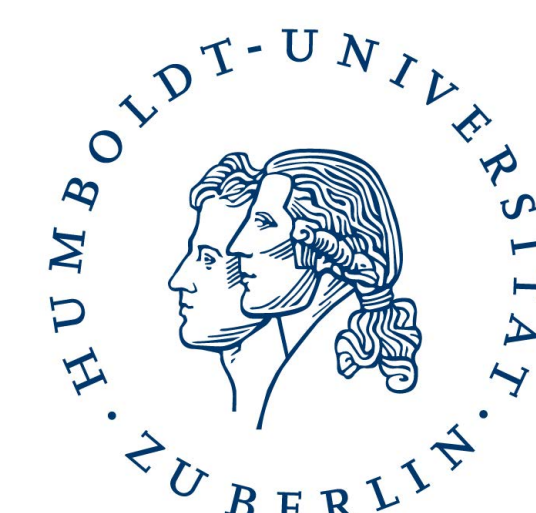


Modellierung linguistischer Forschungsdaten

Was weißt du alles über (d)ein Korpus?



Doktorandentag 2014 Institut für deutsche Sprache und Linguistik 06.10.2014
 Carolin Odebrecht Betreuer: Prof. Dr. Anke Lüdeling, Dr. Laurent Romary

1. Linguistische Forschungsdaten, ihre Metadaten und Modellierung

- Was haben diese Korpora gemeinsam?

Eingrenzung auf historische Korpora aus der Linguistik (Lemnitzer & Zinsmeister 2008, Claridge 2008)

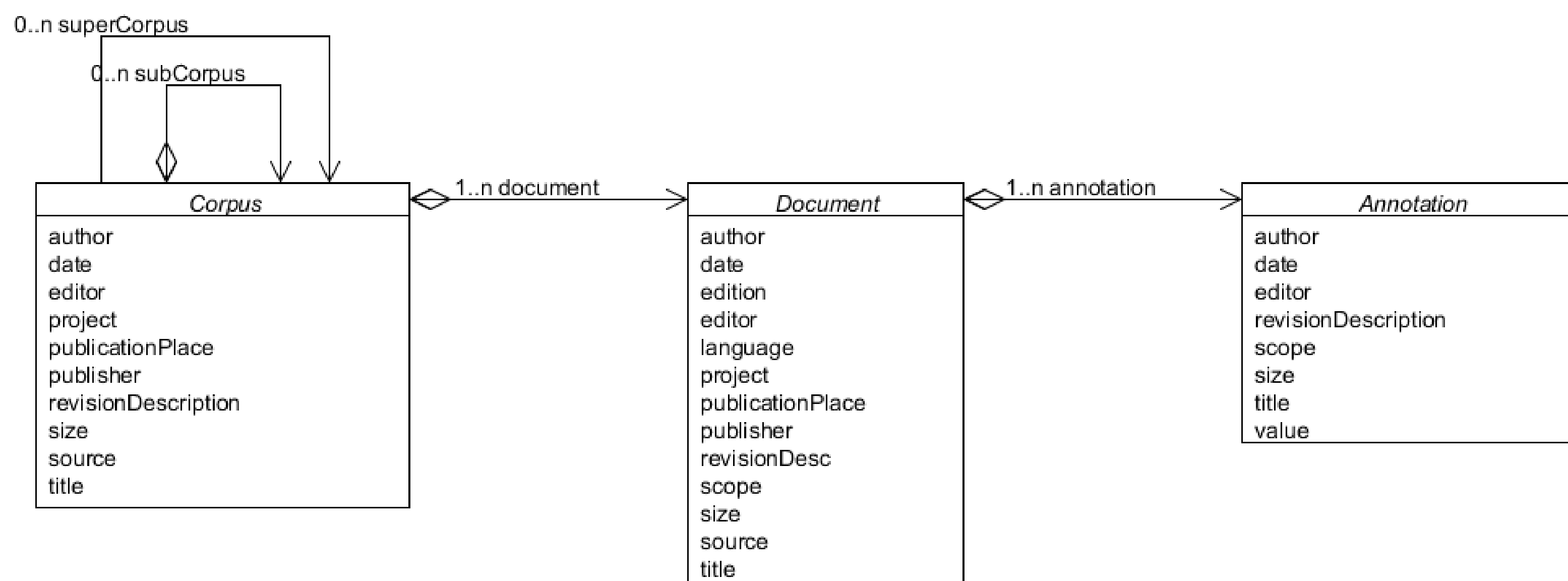
- **Einheiten - Token**
- Annotation
- historische Vorlage für die Digitalisierung

Metadaten linguistischer Forschungsdaten (Odebrecht 2014)

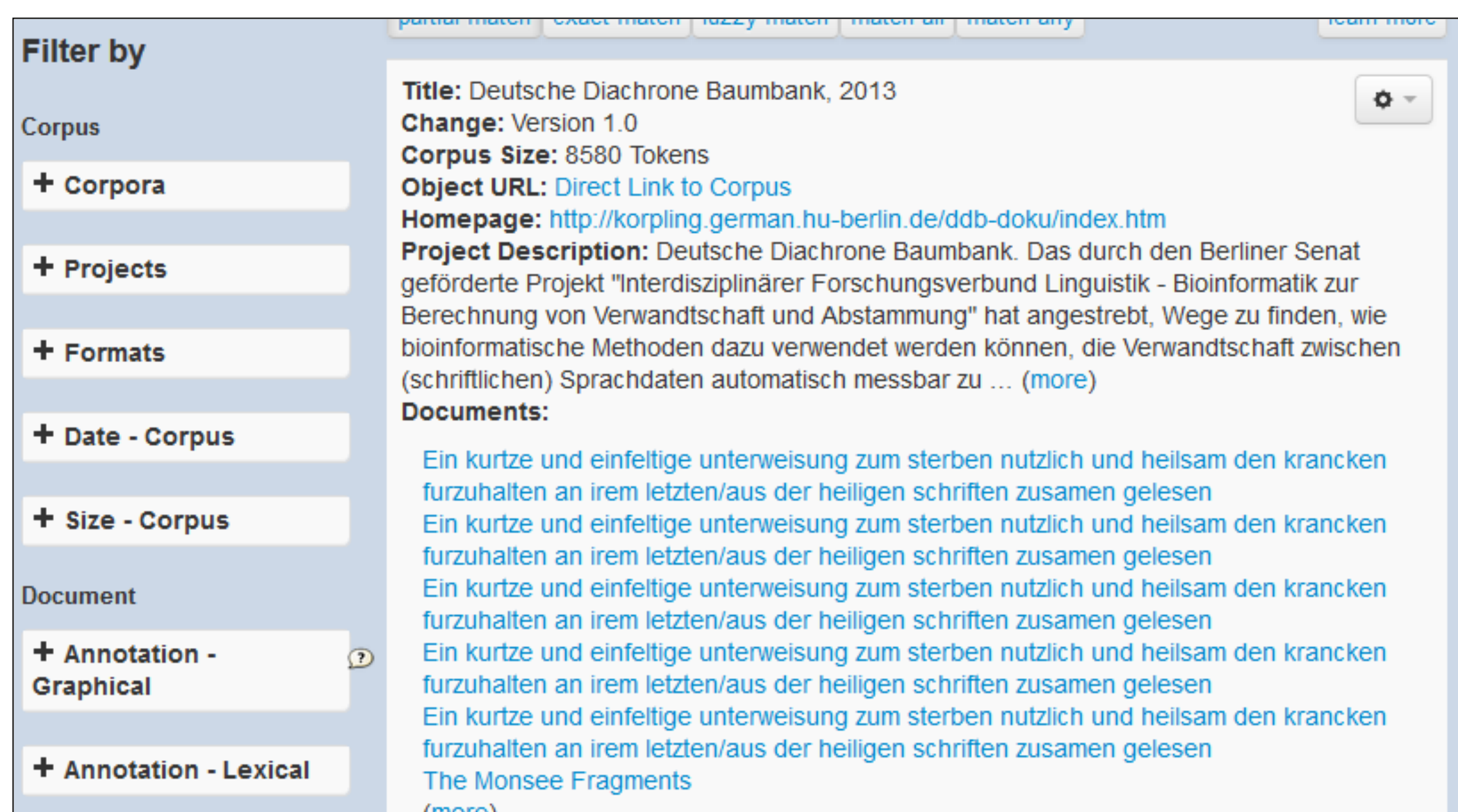
- Daten zum Beschreiben der Daten (Metadaten der Korpora)
- Wofür sollen die Metadaten Korpora beschreiben?
 - Finden von Korpora für eigenen Forschung in einem Repository (Odebrecht et al. 2014)
 - bibliographische Eigenschaften der Texte (Titel, Jahr, Autor, Veröffentlichungsort, Genre etc.)
 - korpuslinguistische Eigenschaften (Format, Projekt, Annotatoren, Annotationsrichtlinien, etc.)
- Wie sollen die Metadaten beschreiben?
 - **einheitliche, strukturierte, umfangreiche Metadaten** historischer Korpora
 - unabhängig von Format
 - unabhängig von Annotation
 - unabhängig von der linguistischen Fragestellung

- Worin unterscheiden sich Korpora?

- Art, Werte, Richtlinien, Tools für Annotationen
- Texte (u.a. Register, Genre, Jahr, Ort)
- Formate (u.a. TIGER XML, EXMARaLDA, Excel, PlainText)
- Segmentierung (einfache oder multiple Segmentierung, vgl. Krause et al. 2012)
- Projekt und Fragestellung (z.B. DDB, HSJ, KAJUK, RIDGES)



2. Anwendung des Modells und technische Umsetzung



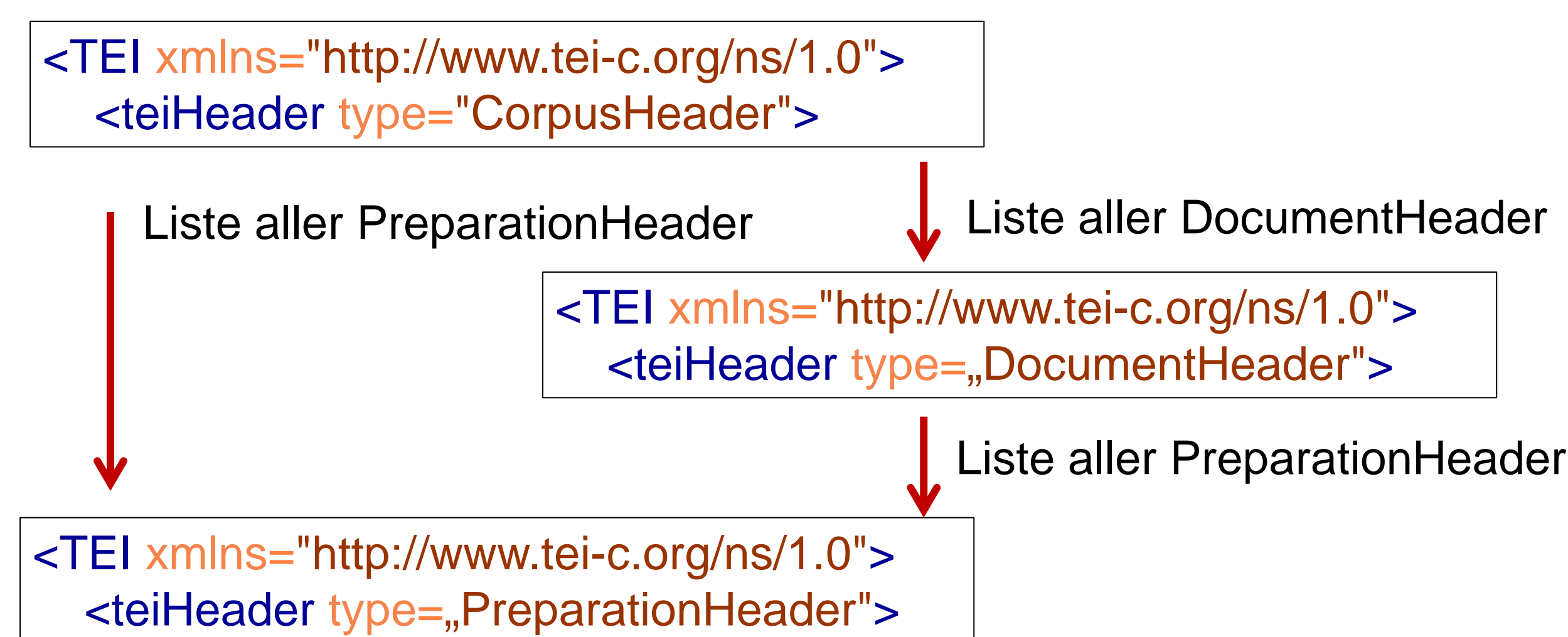
LAUDATIO-Repository

- Facettensuche über Metadaten
- Beschreibung von relevanten Klassen und deren Eigenschaften durch das Metamodell
- Suche und Auswahl eines Korpus nach diesen Eigenschaften
 - u.a. einzelne Facetten für solche Eigenschaften
 - z.B. Annotationsebenen, Größe des Korpus, Projekte
 - alle Eigenschaften sind über die Freitextsuche abrufbar

Umsetzung mit TEI XML

- Annahme: das muss man über sein Korpus wissen – drei Klassen mit ihren Attributen
- jede Klasse erhält einen `teiHeader` (Burnard & Baumann 2008)
 - Corpus – `CorpusHeader`
 - Document – `DocumentHeader`
 - Annotation – `PreparationHeader`
- `teiHeader` sind speziell dafür angepasst (ODD + Schema, vgl. Burnard & Rahtz 2004)
- TEI XML als Format, um diese Metadaten dem Repository zu geben

→ Anzeige der Metadaten für jedes Korpus: **einheitlich, strukturiert und umfassend**
 → Suche nach Korpora über diese Metadaten



3. Neuer Ansatz für die Modellierung von Korpora

Neue Aufgabenstellung

- bislang noch keine Umsetzung von Superkorpus oder Subkorpus im Repository
- noch keine Annotationen, die zwischen Dokumenten bestehen, abgebildet (Anwendung vgl. Berti 2012),
- bislang nur ‚Annotation‘ innerhalb ‚Document‘
- keine Einbindung von ‚Token‘

Klasse ‚Token‘

- **Mengen/Konzeptdefinitionen**
- Annahme: jedes zu beschreibende Korpus besitzt Token
- Token besitzen die Summe aller Attribute der verschiedenen Klassen
- diese Attribute sind auch so referenziert (a für Annotation / d für Document / c für Corpus)
- Token können also über ihre Attribute und deren Werte Gemeinsamkeiten aufweisen

- **Bestimmung der Zugehörigkeit zu einer Klasse: die Menge an gleichen Werten zu einer bestimmten Gruppe an Attributen**

- Relationen via Schnittmengen von Token mit bestimmten Attributen und deren Werten
- Zugehörigkeit zu einer Klasse über gemeinsame Attribute
- Zugehörigkeit zu Instanzen einer Klasse über gemeinsame Werte von Attributen

Token
tID
tsurface
cauthor
cdate
ceditor
cproject
cpublicationPlace
cpublisher
crevisionDescription
csize
csource
ctitle
dauthor
ddate
dedition
deditor
dlanguage
dproject
dpublicationPlace
dpublisher
drevisionDesc
dscope
dsize
dsource
dtile
aauthor
adate
aeditor
arevisionDescription
ascope
asize
atitle
akey
avalue

Welche Eigenschaften besitzt ein Token?

- referenzierbare, kleinste zu annotierende Einheit einer Segmentierung
- ohne Semantik, mit mindestens einer isolierten Eigenschaft
- ...
- ...

Beispiel ‚Token‘ und ‚Document‘ in RIDGES Herbology Corpus

- 154.266 Token in RIDGES Herbology 4.0
- um zu einer Klasse ‚Document‘ zu gehören, müssen Token gleiche Attribute besitzen:
 - dtitle, dauthor und ddate
- um zu einer Instanz der Klasse ‚Document‘ zu gehören, müssen Token diese Gruppe von Attributen mit gleichem Wert besitzen:
 - New Kreuterbuch, Fuchs, 1543

Beispiel ‚Token‘ und ‚Annotation‘ in RIDGES Herbology Corpus

- 154.266 Token in RIDGES Herbology
- um zu einer Klasse ‚Annotation‘ zu gehören, müssen Token folgende gleiche Attribute besitzen:
 - atitle, aauthor und adate
- um zu einer Instanz der Klasse ‚Annotation‘ zu gehören, müssen Token diese Gruppe von Attributen mit gleichem Wert besitzen:
 - pos, Perlitz, 2014-04-01

Nächste Schritte / Ausblick

- tokenbasiertes Modellieren erlaubt einfachere Mengendefinition
- Annotationen, die über die Grenze eines Dokumentes hinaus verweisen, anders zu formulieren: Token, die die Attribute und Werte zu einer Instanz der Klasse ‚Annotation‘ teilen, aber nicht die der Klasse ‚Document‘
- Ausformulieren des Modells, technische Implementation sowie umfangreiches Testen