# Data Modelling for Historical Corpus Annotation

## Cristina Vertan, University of Hamburg

**Textdatenbank und Wörterbuch des Klassischen Maya**

**HerCoRe**

---

Textstruktur-Element, Ebene 4

**Satz1**

ID: TraCES-TestDom-W3

- Textstruktur-Annotation
- Linguistische Annotation
- Edition
- NE

| ' | ə | **m** | m | a | ś | ə | g | a | r | t | u |

Token 1 — **Präp.** — ID: T0>TraCES-TestDom-W3

Token 2 — **Nomen** — ID: T1>TraCES-TestDom-W3

Token 3 — **Pron. Suff.** — ID: T2>TraCES-TestDom-W3

NE (Person, Datum, Ort, Metapher ....) — **Metapher**

⟷ Verlinkung gelöscht     ⟷ Verlinkung eingefügt     ⟷ Verlinkung automatisch generiert

---

## Classical Ethiopic

Gǝˈǝz (a south semitic language) is mentioned for the first time as lnguage of the Aksum Kingdom (approx. 1st-7[th] century. B.C.). The Aksum Kingdom covered the north alpine part of today Ethiopia (Tǝgray) as well as Eritra until the Read Sea.the following centuries, until nowadays Gǝˈǝz was preserved as main language of the Christian Ethiopic Church.

The very rich early Christian literature contains a large number of translations from ancient Geek and Arabic (after the 13[th] century). Gramatical interference phenomena show these influences

**Own Alphabet with syllabic rules**

- Written from left to right
- Complete representation of the vowels (in contrast with Arabic and Hebrew)
- Laryngals and Sibilanta are phonemich and graphemic beliebig interchangeable.

**Noconcatenative Morphologi:** a lexeme is described as a combination of two Elements: the Root and the Pattern (Scheme) .

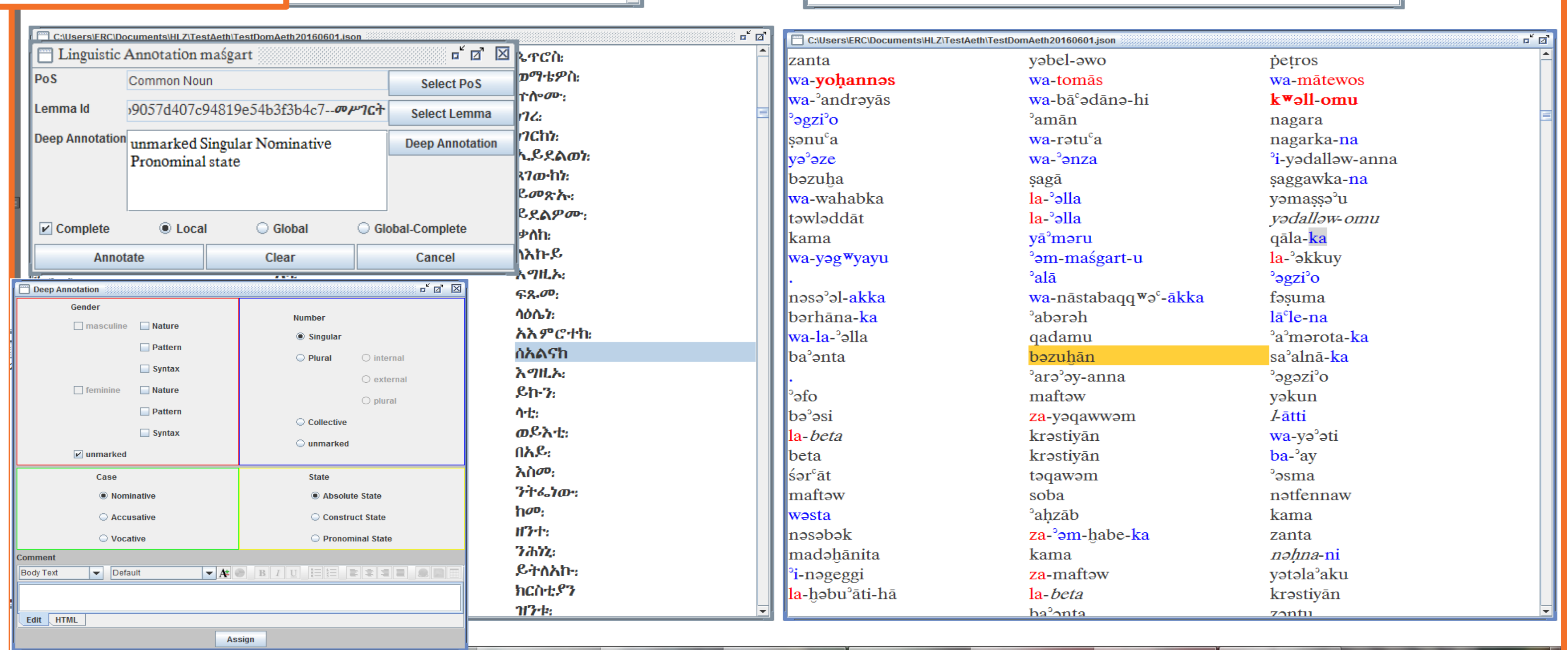$$\sqrt{qds} \times \text{ˈasta}1a22a3a \rightarrow \text{ˈastaqaddasa}$$

---

## Features

- Semi-automatic tokenisation as well as annotation of linguistic features and Nes
- Automatic Sentence recognition
- Visualisation of results of automatic operations
- Post-Annotation correction of the transliterated base Text without lost of annotation information
- Multi-levelAnnotation
- Synchronisation between original and transcribed/transliterated Text

## Implementation

- Java -based client application
- Output-Format JSON
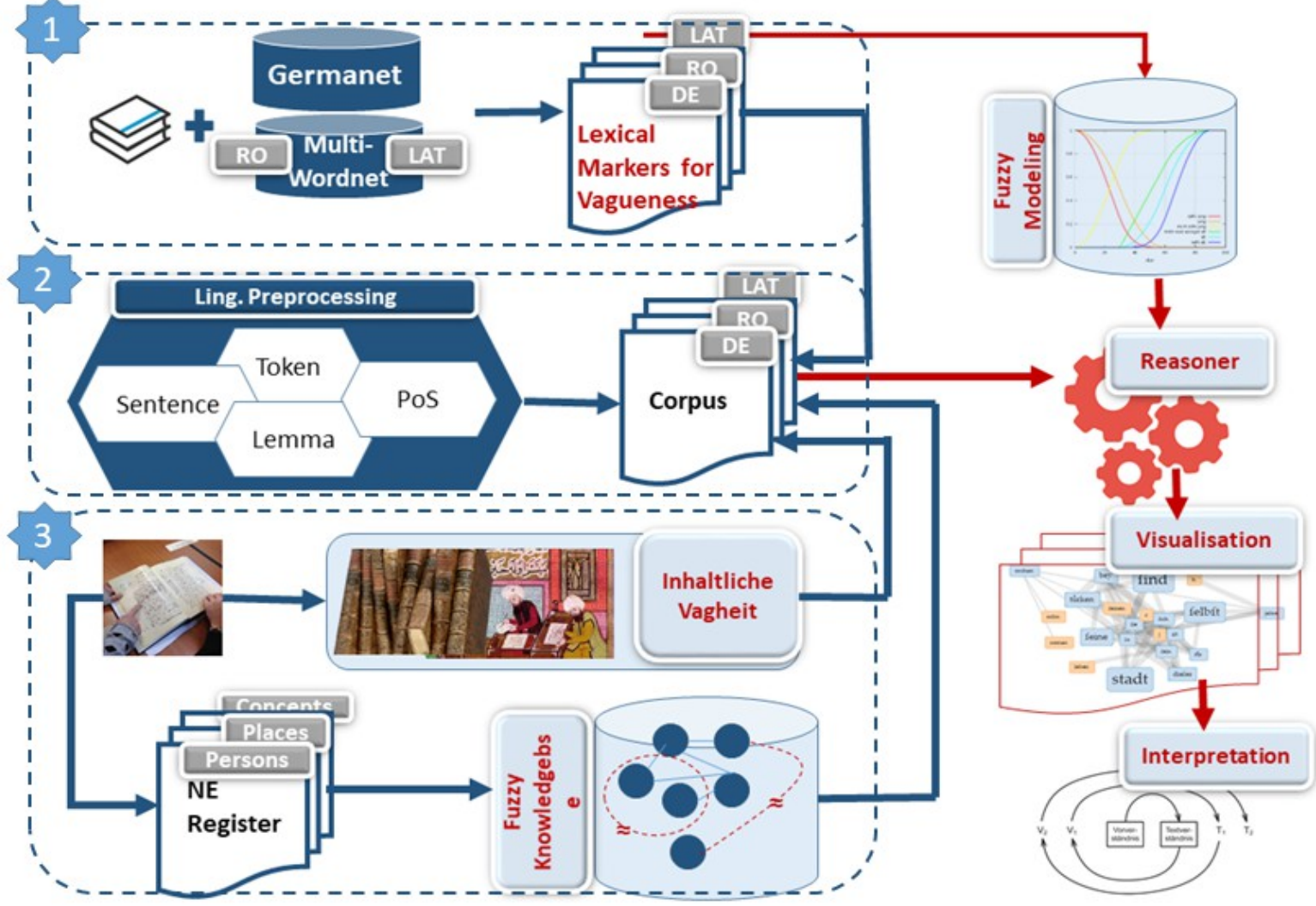- Export TEI, ANNIS, PAULA-XML
- Big-data management

---

## HerCoRe

### Dimitrie Cantemir (1673-1723)

- Prince of Moldavia (historical province) as well as „universal" humanist (linguist, ethnographer, musicologist, historian, writer)
- As member of the Royal Academy in Berlin and at the request of this institution wrote two works :
  - Description of his own country („Descriptio Moldaviae")
  - History of ottoman empire (History of Growth and Decay of Ottoman Empire)
- Original material written in Latin; Both originals were lost already by the end of 18th century
  - Several copies were used as basis for translations into German, English, French, Russian and later in Romanian
  - Sometimes the translation relies on other translation (e.g. first Romanian translation of "Descriptio Moldaviae" was done after the German version from 1774.
- These translations used as reference information about the Ottoman Empire and Romanian provinces until the middle of 19th century, i.e. they give an idea about the reception about this part of the world in Western Europe.

## System Architecture



---

## Word

Id:String
Label_Normalized:String
Label_Original:String
Layout_Information — Layout_Info — Style: Bold/Italic   Script: Fraktur/gr/Lat
Link_to_Entity_Annotation → Entity-Annotation
Links_to_Vagueness_Anotations → Vagueness-Annotation
Link_to_Linguistic_Annotation → Linguistic-Annotation
Links_to_Struture_Annotation → Structure-Anotation

---

### Linguistic Annotation maśgart

PoS: Common Noun — Select PoS
Lemma Id: 9057d407c94819e54b3f8b4c7-- <span>መስገርት</span> — Select Lemma
Deep Annotation: unmarked Singular Nominative Pronominal state — Deep Annotation

☑ Complete   ● Local   ○ Global   ○ Global-Complete

Annotate   Clear   Cancel

### Deep Annotation

**Gender**
☐ masculine — ☐ Nature / ☐ Pattern / ☐ Syntax
☐ feminine — ☐ Nature / ☐ Pattern / ☐ Syntax
☑ unmarked

**Number**
● Singular
○ Plural — ☐ internal / ☐ external / ☐ plural
☐ Collective
☐ unmarked

**Case**
● Nominative
○ Accusative
○ Vocative

**State**
○ Absolute State
○ Construct State
○ Pronominal State

---

## Maya



**9 -Annotations Layers x ?variants /layer**