



# Open Text Collections

Sebastian Nordhoff

Corpus Glosés: de la construction à

l'exploitation automatique

2023-06-28

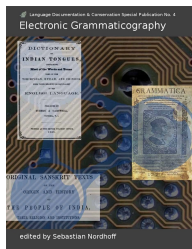


Vol. 2, No. 2 (December 2008), pp. 296-324  
<http://lufic.hawaii.edu/lde/>

## Electronic Reference Grammars for Typology: Challenges and Solutions

Sebastian Nordhoff  
*University of Amsterdam*

Electronic publication offers new possibilities for the creation and exploration of grammatical descriptions. This paper lists values influencing the structure of electronic grammatical descriptions. It then investigates challenges and solutions for a grammar authoring software trying to adhere to these values in the domains of data quality, creation of the description, and exploration of the description. The paper closes by discussing possibilities for the standardization of grammatical descriptions on a macroscopic level, complementing the standardization efforts on a more fine-grained level like GOLD or



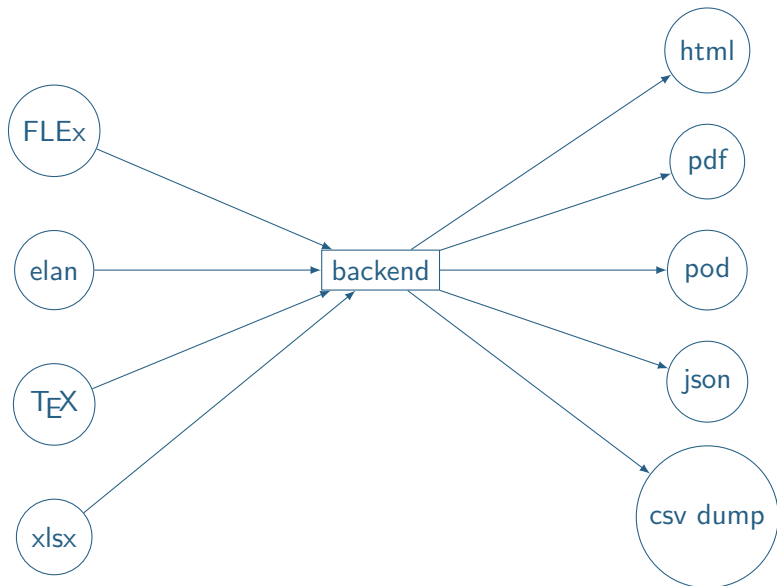
- › work on grammaticography since 2007
- › Language Science Press since 2014
- › starting 2023: open text collections

1. **Dictionaries:** many outlets
2. **Grammatical descriptions:** many outlets
3. **Text collections:** no significant outlets

- › TILA
- › pangloss
- › doreco
- › eopas

- › open
- › prestigious
- › interoperable
- › start 2023-09-15

- › texts
- › edited
- › curated
- › data first
- › community
- › openness
- › prestige







- › CSV for the web (CSVW)
- › text-based
- › simple
- › expandable
- › easy to version

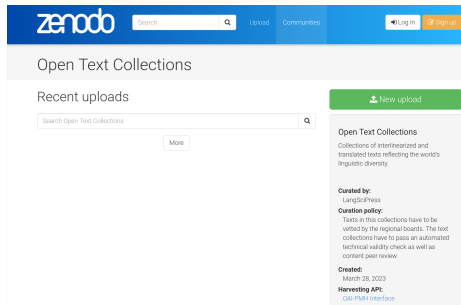
	A	B	C	D
1	Tire la bobinette, la chevillette cher-ra	pull.IMP DEF.F.SG bobbin DEF.F.SG latch fall-FUT.3SG	Pull the bobbin and the latch will fall.	stan1290
2				

- › three obligatory columns
  - › vernacular
  - › glosses
  - › translation
- › can be expanded by further columns as necessary
- › Leipzig Glossing Rules
- › correspondences and constraints between cells in the same row

---

```
> dump (csv, json-ld, nq, rdf)
> pdf
    > printed pdf = book
> html
> query interface
    > see https://imtvault.org/?q=coconut
```

- › preparation on GitHub
- › DOI via the  
GitHub-Zenodo bridge
- › review via GitHub  
issues
- › collections approved by  
the regional boards get  
a new release (1.0),  
archived on Zenodo,  
and are accepted in  
the relevant Zenodo  
communities.

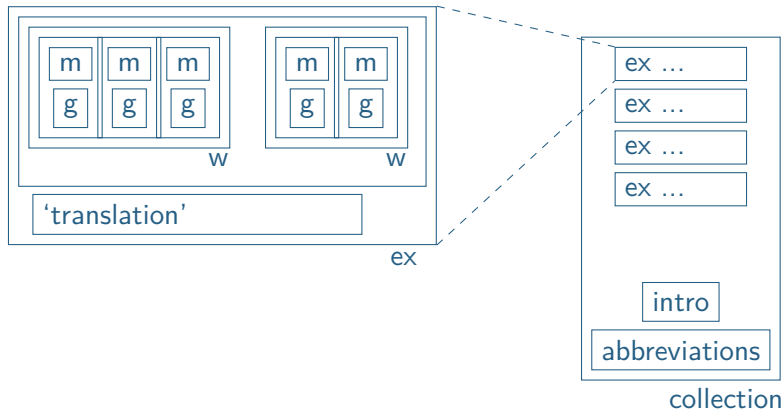


## } technical review

- } number of words match between vernacular and glosses
- } number and type of morphemes match between vernacular and glosses
- } abbreviations are either Leipzig Glossing Rules or listed in separate document

## } content review

- } regional boards with area specialists
- } precise setup to be debated



- › **findable**: registered, linked, metadata
- › **accessible**: no paywalls
- › **interoperable**: many different formats
- › **reusable**: open license

- › 3 year grant from DFG, 2023-2026, 1.5 FTE
- › after that consortial funding via Language Science Press



- › Mandana
- › Christian
- › Sebastian
- › student assistants

- › **Oceania** Christian Döhler, Kilu von Prince (Düsseldorf)
- › **Africa** Alena Witzlack-Makarevich (Jerusalem), Jeff Good (Buffalo)
- › **Eurasia** Michael Rießler (Joensuu)
- › **South America** Matt Coler (Groningen), Nick Emle (Groningen)
- › **Caucasus** Diana Forker (Jena)

## Planned text collections

---

- |                           |                 |
|---------------------------|-----------------|
| › Komnzo                  | › Hinuq         |
| › Bine                    | › Sanzhi        |
| › Daakaka                 | › Chirag Dargwa |
| › Dalkalaen               | › Tabasaran     |
| › Muylaque Aymara         | › Gawarbatl     |
| › Iquito                  | › Palula        |
| › two Amazonian languages | › Saek          |
| › Kawesqar                |                 |