# Using Linked Data Techniques for Creating an isiXhosa Lexical Resource - a Collaborative Approach

Thomas Eckart
Dirk Goldhahn
*Natural Language Processing Group*
*University of Leipzig, Germany*

Bettina Klimek
*AKSW Group*
*University of Leipzig, Germany*

Sonja Bosch
*Department of African Languages*
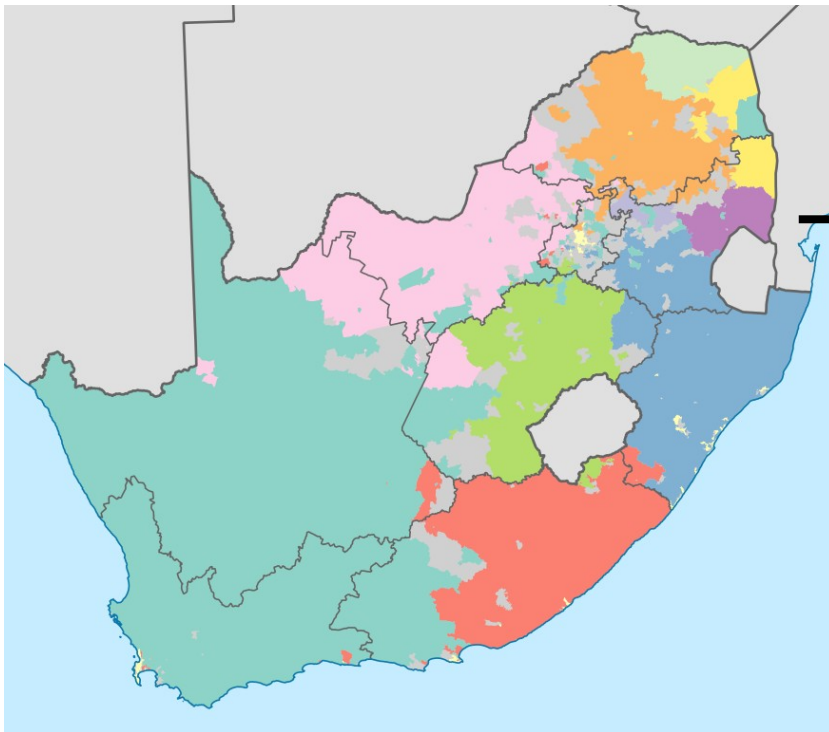*University of South Africa*

# Agenda

- Language Diversity in South Africa
- isiXhosa Data Set
- Bantu Language Model (BLM)
- CLARIN integration and Dissemination
- Further work

# Language Diversity in South Africa

- 11 official languages, including 9 Bantu languages



**Source**: Stats SA

# isiXhosa in South Africa

- Spoken in the Eastern Cape and Western Cape regions

- Approximately 8.1 million first language speakers (16% of South Africa's population)

- Nguni language: conjunctive orthography where bound morphemes are attached to words

- Resource scarce language, especially regarding lexical resources; no machine-readable lexicons freely available

# isiXhosa Data Set (1/2)

- isiXhosa dictionary
  - Based on data by J.A. Louw and S. Bosch
  - Digitization starting with index cards
  - Still work in progress, finished until Mid 2019
  - Available data:
    - Part of speech, noun classes and prefixes, singular/plural forms, English translations
  - Current status:
    - ~4,000 nouns, ~2,700 verbs, ~8,000 translations
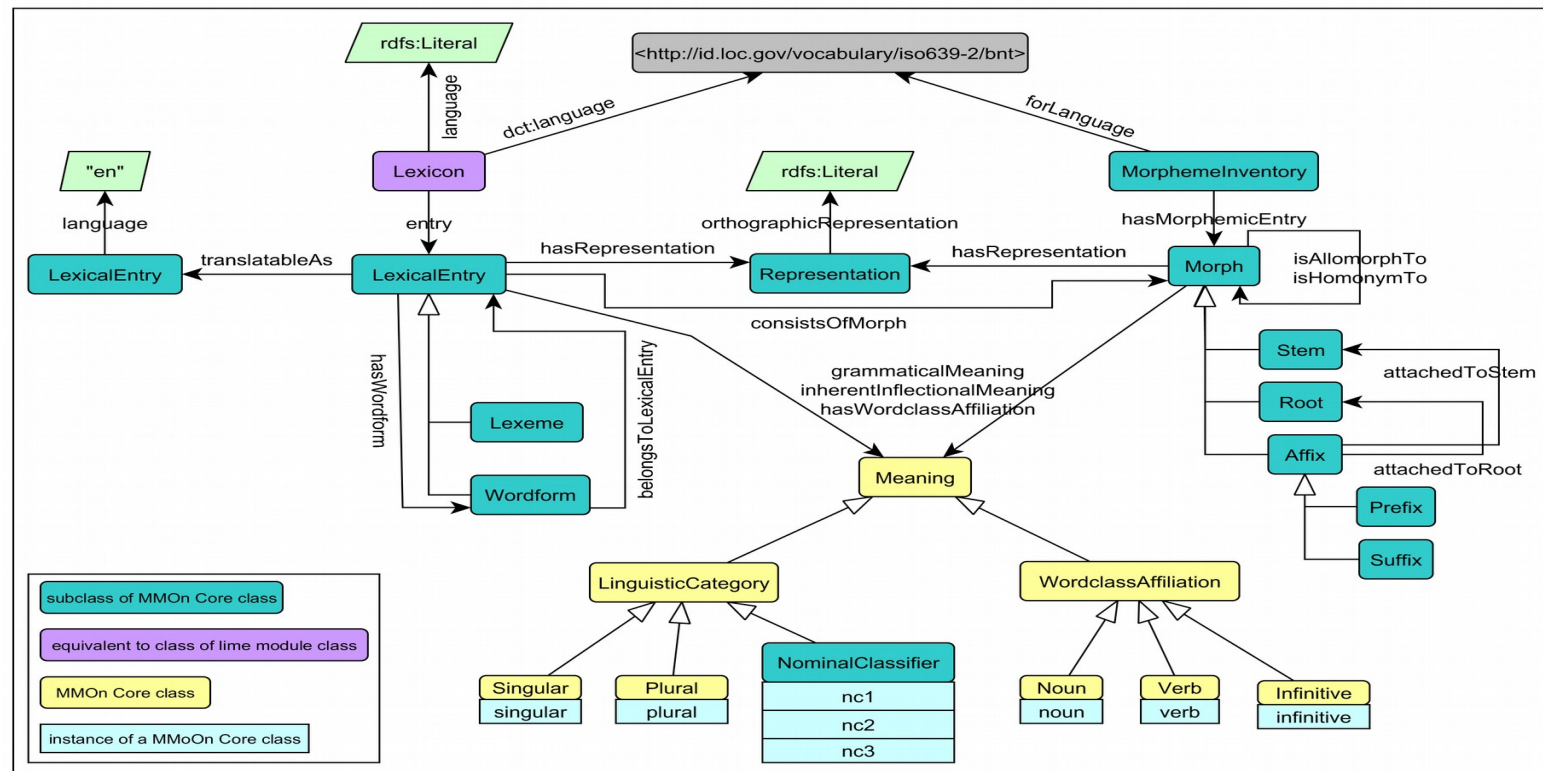
# isiXhosa Data Set (2/2)

- isiXhosa nouns (excerpt):

| | A | B | C | D | E | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|
| 1 | abelo | noun | is | 7 | iz | 8 | | lot | share |
| 2 | abhalala | noun | is | 7 | iz | 8 | | striped Sotho blanket | |
| 3 | abhobho | noun | is | 7 | iz | 8 | | chasm | yawning cavity |
| 4 | abi | noun | um | 1 | ab | 2 | | divider | distributor |
| 5 | abonkolo | noun | is | 7 | iz | 8 | | tadpole | |
| 6 | aci | noun | is | 7 | iz | 8 | | motto | pet saying |
| 7 | adlunge | noun | is | 7 | iz | 8 | | wagon tree | Protea grandiflora |
| 8 | afobe | noun | is | 7 | iz | 8 | | idioms | |
| 9 | akhi | noun | is | 7 | iz | 8 | | formative | morpheme |
| 10 | akhiwo | noun | is | 7 | iz | 8 | | building | |
| 11 | akhiwo | noun | is | 7 | iz | 8 | | structure | building |
| 12 | akhiwo sesikhumbuzo | noun | is | 7 | iz | 8 | | monument | |
| 13 | akhono | noun | is | 7 | iz | 8 | | dexterity | ability to manipulate |
| 14 | alam | noun | is | 7 | iz | 8 | | poor person | |
| 15 | alathandawo | noun | is | 7 | iz | 8 | | agentative | demonstrative pronoun |
| 16 | alathisi | noun | is | 7 | iz | 8 | | demonstrative | |
| 17 | alelo | noun | is | 7 | iz | 8 | | prohibition | opposition |
| 18 | ambalo | noun | is | 7 | iz | 8 | | necklace | |
| 19 | ambantlanya | noun | is | 7 | | | | outcry | gist |

# Bantu Language Model (1/2)

- Structured using Bantu Language Model (BLM)
  - RDF/OWL-based ontology
  - Based on the *Multilingual Morpheme Ontology* (MMoOn)

**Model for lexical data of Bantu languages.**

# Bantu Language Model (2/2)

- Access via dedicated Web page, SPARQL endpoint and downloads

**Source: https://rdf.corpora.uni-leipzig.de/lexeme_ibhatata_n.html**

# CLARIN Integration and Dissemination

- Joint work of three partners, providing their linguistic and technical expertise

- Resource currently hosted at CLARIN-D centre Leipzig and GitHub
  - https://rdf.corpora.uni-leipzig.de
  - https://github.com/MMoOn-Project/OpenBantu

- "Final" resource hosted at *SADiLaR Resource Catalogue*

# CLARIN Integration and Dissemination

- Usage of Handles & CMDI metadata
- Strategy for integrating LLOD in CLARIN?

# Dissemination



Model for lexical data of Bantu languages.

Lexeme: **bhatata**

▲ Morphosyntax

**Word class:** noun
**Inflected forms:** amabhatata, ibhatata, iibhatata, oobhatata, ubhatata

▲ Translations

• sweet potato *(English)*
• sweet potato plant Ipomoea batatas *(English)*

▲ Examples:

• Kodwa, phambi kokuba le ndoda ingcamlise nabani na, yayicela umntu lowo ukuba ayiphe ibhatata, ilungwana lenyama nentwana yetyuwa kunye nepepile. (nalibali.mobi, *crawled on 26/07/2017*)

▲ Records in other dictionaries with similar meaning

• hlaza *(isiNdebele)*
• bambayila *(isiNdebele)*

**XhosaDict**

Wort: **bhatata** (noun)

**Flexionsformen:**
• amabhatata
• ibhatata
• iibhatata
• oobhatata
• ubhatata

**Übersetzungen:**
• sweet potato
• sweet potato plant Ipomoea

**Beispiele:**
• [...] ukuba ayiphe ibhatata, ilungwana [...]

# Further Work

- Continuing work on the data and QA

- Combination with open text corpora

- Open data services for a federated lexicographical infrastructure

  – Prototype: CBOLD data for lexical alignment based on English translations

- Annotation services: WebLicht?

- ?

# Thank you!

Questions? Remarks?