

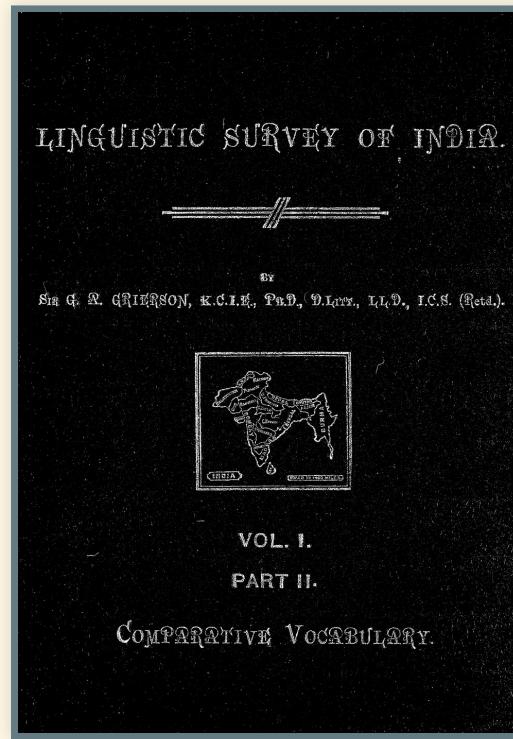
Linguistic Survey of India and Polyglotta Africana: Two Retrostandardized Digital Editions of Large Historical Collections of Multilingual Wordlists

Robert Forkel, Johann-Mattis List, Christoph
Rzymski, Guillaume Segerer

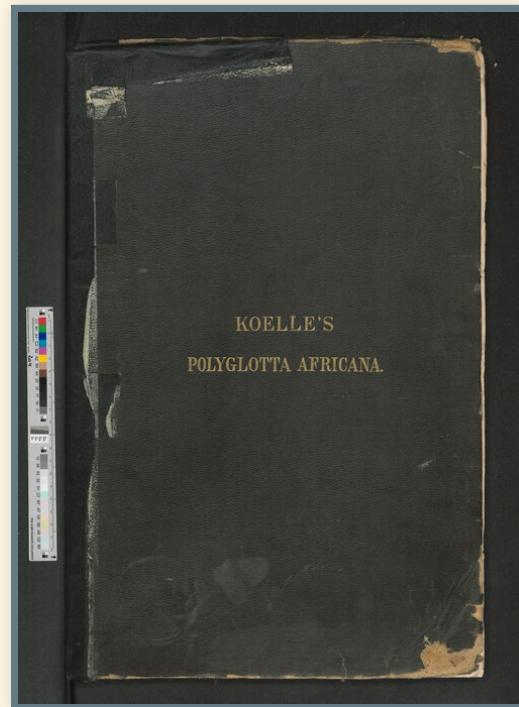
Historical Data for Historical Linguistics

- comparative wordlist have been collected for more than 230 years
- importance of comparative lexical data has increased with automated analysis methods like cognate detection or phylogenetic methods
- automated re-use requires standardized data

The Linguistic Survey of India



Polyglotta Africana



CLDF - an interoperable data format

- tabular data in CSV files
- typed data thanks to schema description according to CSVW
- domain-specific semantics via an ontology
- shared domain data via reference catalogs

Language Identification

- Glottolog - the reference catalog for languoids - is referenced via
- Glottocodes specified using the CLDF glottocode property in a
- CLDF LanguageTable component.

Concept Mapping

- Concepticon - the reference catalog for concepts - is referenced via conceptset IDs specified using the CLDF concepticonReference property in a CLDF ParameterTable component.

Orthography Conversion

- Original orthography is converted via orthography profiles to
- segmented phonetic transcriptions linked to
- CLTS - the reference catalog for transcription systems.

CLDF Creation

- “Raw” digital data is fed into
- `cldfbench` where it is enriched with
 - language and concept mappings and segmented via
 - orthography profiles
- resulting in CLDF datasets.

Orthography Profiles

The screenshot shows a GitHub repository interface. On the left, the 'Files' sidebar displays the repository structure under the 'master' branch. The 'etc/orthography' directory contains four files: ABOR.tsv, AHI.tsv, AHOM.tsv, and AIMOL.tsv. The 'AHI.tsv' file is currently selected and shown in the main pane.

The main pane title is [lsi / etc / orthography / AHI.tsv](#). It includes a 'LinguList' update information button and tabs for Preview, Code, and Blame. The Preview tab shows 84 lines (84 loc) and 823 Bytes. A search bar is available to search the file content.

The AHI.tsv file content is displayed as a table:

	Grapheme	IPA	Frequency
2	^ń	ŋ	2
3	o	o	71
4	°	NULL	135
5	\$	NULL	16
6	^ch'	tʃ̚	2

Validation

The screenshot shows the GitHub Actions page for the repository `lexibank / lsi`. The left sidebar lists actions like "All workflows", "Workflows", "CLDF-validation" (which is selected), and "Management". The main area displays the `CLDF-validation` workflow with its configuration file `cldf-validation.yml`. It shows 8 workflow runs, with two recent ones highlighted:

- ready for release**: CLDF-validation #8: Commit [de3b554](#) (master) pushed by xrotwang, 7 months ago, 59s ago.
- add coordinates for dialects...**: CLDF-validation #7: Commit [ac29409](#) (master) pushed by xrotwang, 7 months ago, 1m 12s ago.

Visualization: HTML

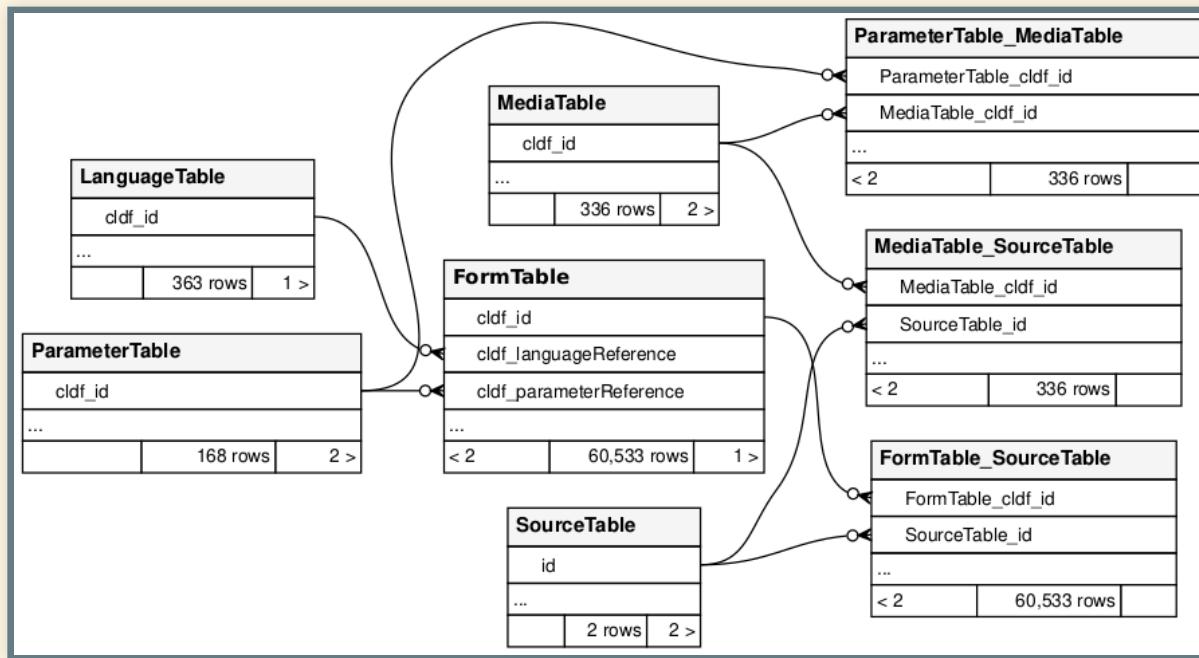
Wordlist CLDF dataset derived from Grierson's "Linguistic Survey of India" from 1928

CLDF Metadata: [cldf-metadata.json](#)

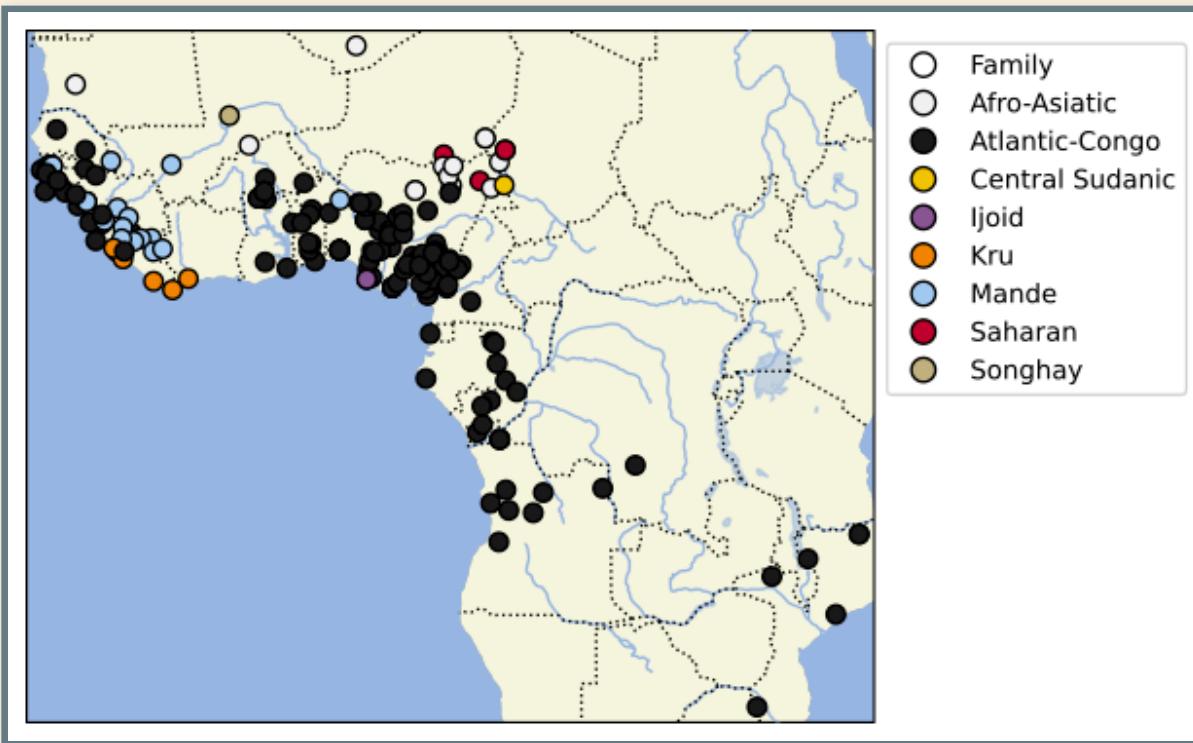
Sources: [sources.bib](#)

property	value
dc:bibliographicCitation	Grierson, George Abraham (1928): Linguistic Survey of India. Comparative Vocabulary. Calcutta: Government of India Central Publication Branch.
dc:conformsTo	CLDF Wordlist
dc:format	1. http://concepticon.clld.org/contributions/Grierson-1928-168
dc:identifier	https://lsi.clld.org
dc:license	https://creativecommons.org/licenses/by/4.0/
dcat:accessURL	https://github.com/lexibank/lsi
prov:wasDerivedFrom	<ol style="list-style-type: none">1. lexibank/lsi ac294092. Glottolog v4.83. Concepticon v3.1.04. CLTS v2.2.0

Visualization: ERD



Visualization: Coverage map



Visualization: clld

Koelle's Polyglotta Africana -- Chromium

Koelle's Polyglotta Africa +

polyglottafricana.clld.org/valuesets/I-A-1-1_one

POLYGLOTTA AFRICANA.

Words for One in I-A-1 Fúlup

Source Broad IPA

	Source	Broad IPA
fánqd	fanɔ̄d	

	English.	One.	Two.	Three.	English.	One.	Two.	Three.
1. A. 1. Fulup	fánqd	fúgaptén & fáten	fúfontén	k. Ifé	éni	édzi	éta	
2. Filham	ánod	kágwá	kifégi	l. Ondó	íng	édzi	éta	
B. 1. Béla	puljó	kétaw	kéhawas	m. Ðékirí	ménę	méds	méta	
				z. Tesla	ínwe	z. i.	z.	

<https://pic.sub.uni-hamburg.de/kitodo/PPN862704383/00000042.tif>

Visualization: IPA chart

Aggregation: CLDF SQL

```
1 ATTACH DATABASE "phoible.sqlite" AS phoible;
2 ATTACH DATABASE "clts.sqlite" AS clts;
3 ATTACH DATABASE "lsi.sqlite" AS lsi;
```

```
5 CREATE TEMP VIEW lsigraphemes AS
6 SELECT
7     DISTINCT grapheme
8 FROM
9 (
10    WITH split(grapheme, segments) AS (
11        SELECT
12            '',
13            f.cldf_segments || ''
14        FROM lsi.formtable AS f, lsi.languagetable AS l
15        WHERE f.cldf_languagereference = l.cldf_id AND l.cldf_glottocode = 'mala1464'
16    UNION ALL SELECT
17        substr(segments, 0, instr(segments, ' ')),
18        substr(segments, instr(segments, ' ') + 1)
19    FROM split
20    WHERE segments != ''
21    ) SELECT grapheme FROM split
22    WHERE grapheme != ''
23 );
```

```
25 CREATE TEMP VIEW phoiblegraphemes AS
26 SELECT
27     DISTINCT c.cltsgrapheme AS grapheme
28 FROM
29 (
30     SELECT v.cldf_value AS grapheme
31     FROM phoible.valuetable AS v
32     WHERE cldf_languagerefERENCE = 'mala1464' AND contribution_id = 1762
33 ) AS p
34 JOIN
35 (
36     SELECT phoible.grapheme AS phoiblegrapheme, clts.grapheme AS cltsgrapheme
37     FROM clts."data/graphemes.tsv" AS phoible, clts."data/sounds.tsv" AS clts
38     WHERE phoible.dataset = 'phoible' AND phoible.name = clts.name
39 ) AS c
40 ON c.phoiblegrapheme = p.grapheme;
```

```
42 SELECT lsi.grapheme, 'LSI', clts.name
43 FROM lsigraphemes AS lsi
44 JOIN
45     clts."data/sounds.tsv" AS clts
46 ON clts.grapheme = lsi.grapheme
47 WHERE clts.name LIKE '%vowel' AND lsi.grapheme NOT IN phoiblegraphemes;
48
49 SELECT phoible.grapheme, 'PHOIBLE', clts.name
50 FROM phoiblegraphemes AS phoible
51 JOIN
52     clts."data/sounds.tsv" AS clts
53 ON clts.grapheme = phoible.grapheme
54 WHERE clts.name LIKE '%vowel' AND phoible.grapheme NOT IN lsigraphemes;
```

Glyph	Dataset	CLTS
ʌ	LSI	unrounded_open_mid_back_vowel
ʌ:	LSI	long_unrounded_open_mid_back_vowel
a	PHOIBLE	unrounded_open_front_vowel
a:	PHOIBLE	long_unrounded_open_front_vowel
æ	PHOIBLE	unrounded_near_open_front_vowel
ɨ	PHOIBLE	unrounded_close_central_vowel
ʊ	PHOIBLE	rounded_near_close_near_back_vowel

Analysis

- LingPy can read CLDF and compute cognates or alignments
- BEASTling can read cognate coded CLDF to compute phylogenetic trees

Conclusion

- We have established a workflow to retro-digitize legacy wordlists.
- Thereby providing another angle from which to attack the data scarcity problem in computational linguistics.

