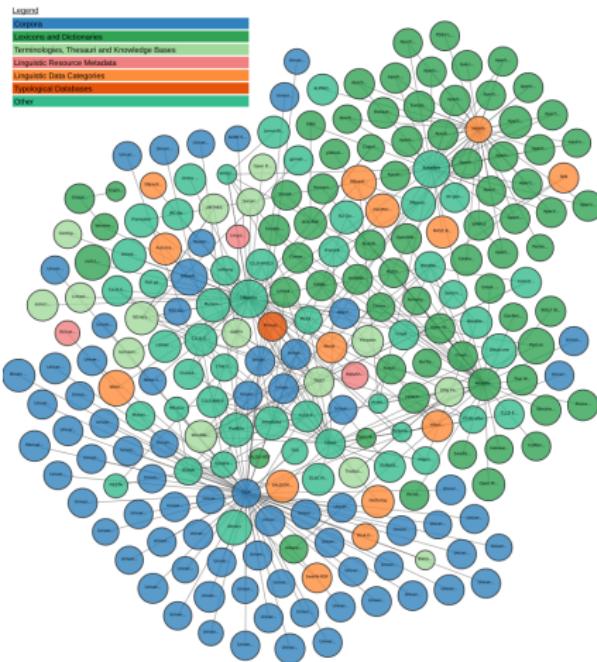




IMTVault: Linked Data from Open Access L^AT_EX books

SN & TK
June 2, 2022

The LLOD cloud



The Linguistic Linked Open Data Cloud from llo-d.cloud.net



structured language data
for NLP,
WordNet,
Thesauri,
Dictionaries,
Treebanks,
Annotated corpora

For how many
languages?

Coverage of NLP

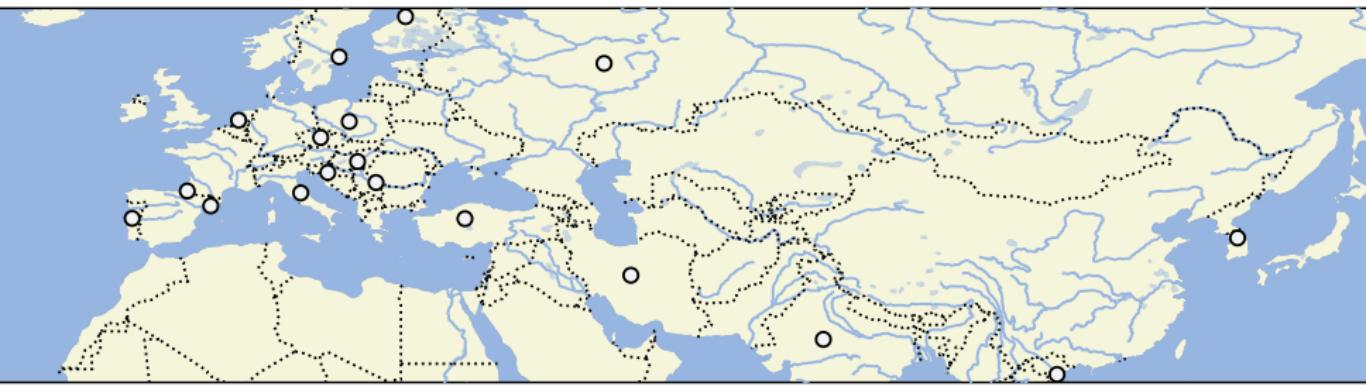
Joshi analyzed the languages covered in NLP research and resources.

Class	criteria			example	# lgs	%
	unlabeled data	labeled data				
5	winners	good	good	Spanish	7	0.28
4	underdogs	good	insufficient	Russian	18	1.07
3	rising stars	good	none	Indonesian	28	4.42
2	hopefuls	?	smallish sets	Zulu	19	0.36
1	scraping-bys	smallish	none	Fijian	222	5.49
0	left-behinds	none	none	Warlpiri	2191	88.38
-1	not included	?	?	Komnzo	6404	-

Joshi's group 5: winners



Joshi's group 4: underdogs



Joshi's group 3: rising stars



Joshi's group 2: hopefuls



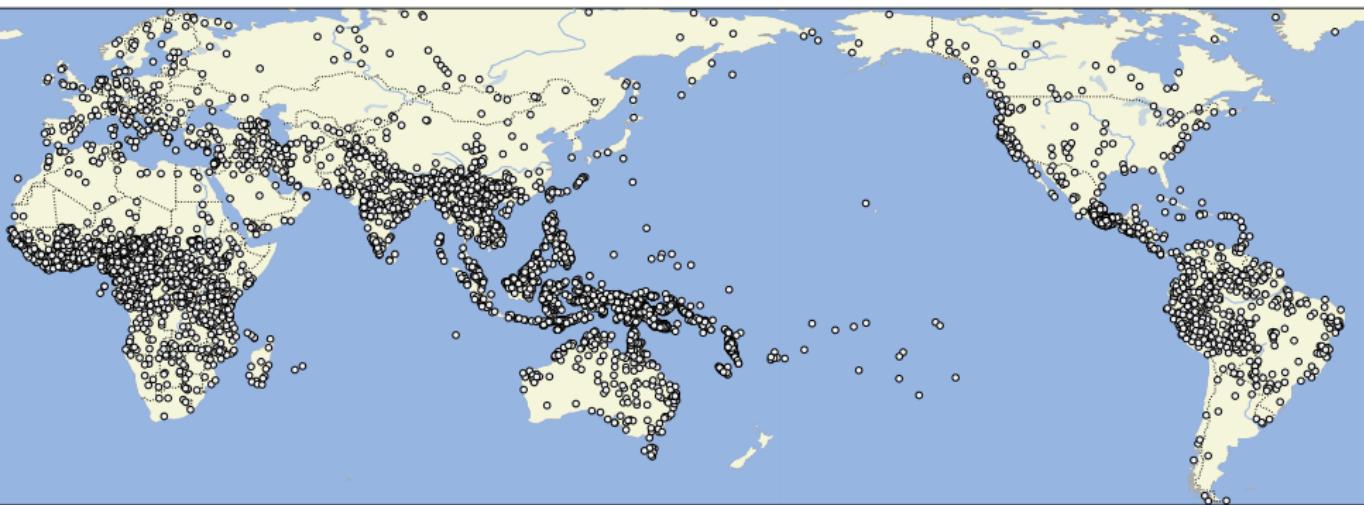
Joshi's group 1: scraping-bys



Joshi's group 0: left-behinds



Group -1: there is such a language??



What do we have instead of “classic” NLP stuff?

- › For many of the languages of groups 0 and –1, we do indeed have some resources
 - › typological databases (WALS, APICS)
 - › Endangered Language Archives (ELAR, PARADISEC, TLA, AILLA)
 - › scientific publications (books, articles)

Common datastructure: IGT

› IGT: interlinear glossed text

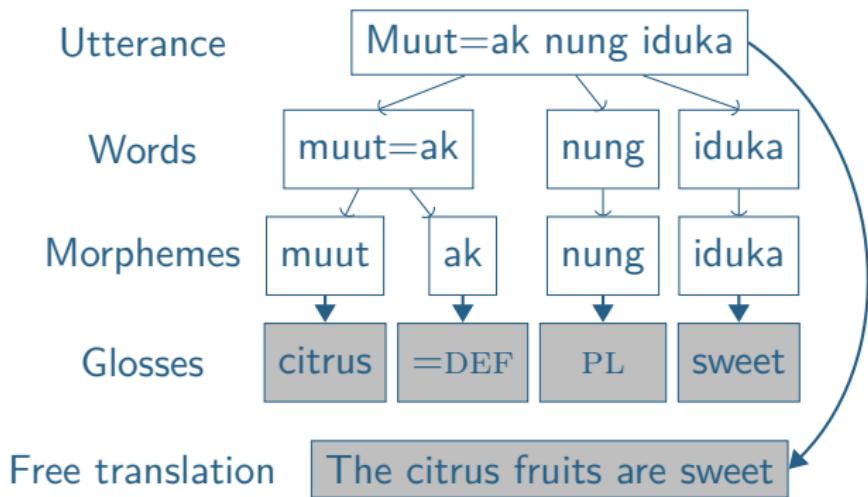
(1) KamangSchapper, fieldnotes

Muut=ak nung iduka.

citrus=DEF PL sweet

'The citrus fruits are sweet.'

Common datastructure: IGT





THE ATLAS OF PIDGIN AND CREOLE LANGUAGE STRUCTURES ONLINE

[Home](#)[Languages](#)[Features](#)[WALS-APiCS](#)[Surveys](#)[Examples](#)[Sources](#)

Example 22-78

Bai mitupela i ringim taksi.

Bai mitupela i ring-im taksi.

FUT 1DU.EXCL PM ring-TR taxi

'We'll ring a taxi.'

Type:

naturalistic spoken

Source:

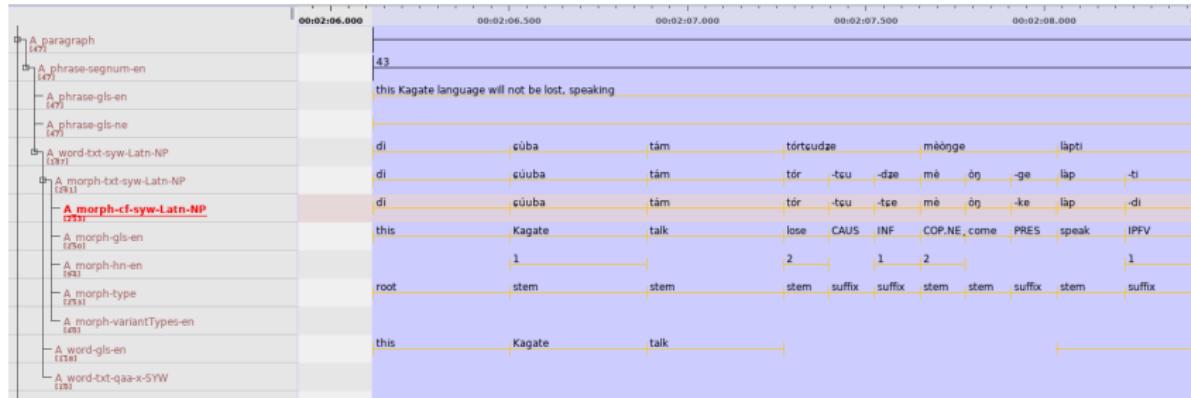
Dutton and Thomas 1985: 88

Nordhoff & Krämer

(Chiarcoslonov)

14/30

ELAN files from endangered language archive



This talk: scientific publications

```
\langinfo{Kamang}{}{Schapper, fieldnotes} \\
\gll Muut=ak nung iduka. \\
    citrus=\textsc{def} \textsc{pl} sweet \\
\glt 'The citrus fruits are sweet.'
```

-
- › open access publisher in linguistics founded 2014
 - › 180 published books
 - › all books licensed as CC-BY
 - › tex codes for all books available no GitHub
 - › \LaTeX package gb4e used for interlinear glossed text

A grammar of Gyeli

A grammar of Gyeli

Nadine Grimm

Synopsis

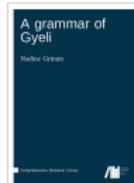
This grammar offers a grammatical description of the Agélo variety of Gyeli, an endangered Bantu (ABl) language spoken by 4,000–5,000 “Pygmy” hunter-gatherers in southern Cameroon. It represents one of the most comprehensive descriptions of a northwestern Bantu language.

The grammatical description, which is couched in a form-to-function approach, covers all levels of language, ranging from Gyeli phonology to its information structure and complex clauses.

It draws on nineteen months of fieldwork carried out as part of the ‘*Bayesi/Bakola*’ OdBeL (Documentation of Endangered Languages) project between 2010 and 2014. The resulting multimodal corpus from that project, which includes texts of diverse genres such as traditional stories, narratives, multi-party conversations and dialogues, procedural texts, and songs, provides the empirical basis for the grammatical description. The documentary text collection, supplemented by data from elicitation work, questionnaires, and experiments, are accessible in the [‘*Bayesi/Bakola* collection’ of the Language Archive](#). With additional ethnographic, sociolinguistic, diachronic, and comparative remarks, the grammar may appeal to a wider audience in general linguistics, typology, Bantu studies, and anthropology.

In 2019, the grammar received the Pâline Award by the Association for Linguistic Typology.

Statistics



[PDF ↗](#)

[Bibliography ↗](#)

[Buy from amazon.de ↗](#)

[Buy from amazon.co.uk ↗](#)

[Buy from amazon.com ↗](#)

[Collaborative reading on
Papertree ↗](#)

[LaTeX source on GitHub](#) →

langsci / 298 · Public

< Code Issues Pull requests Actions Projects Wiki

main · 298 / chapters / CH7.tex

Glottotopia copy from overleaf

1 contributor

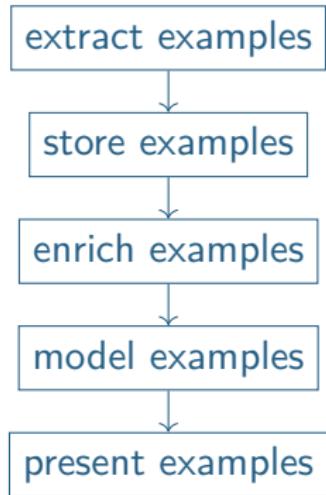
2871 lines (2042 sloc) | 186 KB

```

1 \chapter{Simple clauses}
2 \label{sec:SC}
3 
4 
5 
6 In this chapter, I describe the different types of simple clauses in Gyeli.
7 
8 
9 
10 
11 
12 \section[Copula constructions]{Non-verbal and verbal copula constructions}
13 \label{sec:nonverbcl}
14 
15 Gyeli has copula clauses with both non-verbal and verbal copula construction
16 
17 \textsf{\textless}ea\label{John}\textsf{\textgreater}
18 \gll John (\textsf{\textless}bfrseries ni) a-kubwa \\
19 \$emptyset\$1 (\textsf{\textless}PN\rangle) (\textsf{\textless}NP\rangle) 3-big \\
20 \trans 'John is big'
21 \z

```

pipeline



extract examples

KamangSchapper, fieldnotes

Muut=ak nung iduka.

citrus=DEF PL sweet

'The citrus fruits are sweet.'

```
\langinfo{Kamang}{}{Schapper, fieldnotes}
```

```
\gll Muut=ak nung iduka.
```

```
citrus=\textsc{def} \textsc{pl} sweet
```

```
\glt `The citrus fruits are sweet.'
```

Store examples

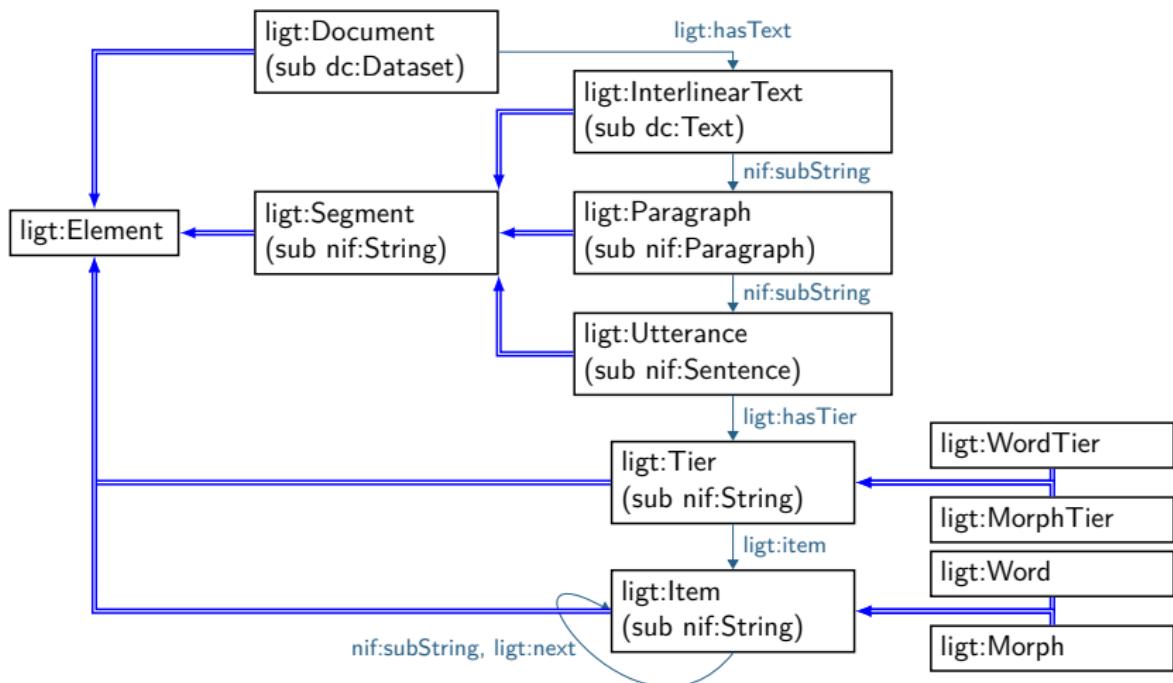
```
{"ID": "w109-cb9806ea53",
"book_ID": 157,
"book_URL": "https://langsci-press.org/catalog/book/157",
"book_metalanguage": "eng",
"book_title": "The Alor-Pantar languages2",
"categories": ["def", "pl"],
"imtwordsbare": [
    "citrus=DEF",
    "PL",
    "sweet"
],
"label": "Muut=ak nung iduka.",
"srcwordsbare": [
    "Muut=ak",
    "nung",
    "iduka."
],
"trs": "The citrus fruits are sweet."
}
```

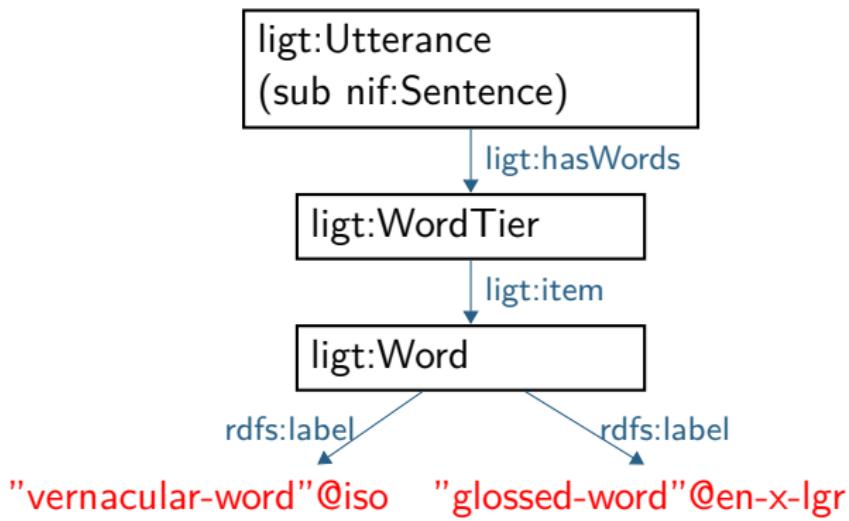
Enrich examples

- › Use Glottolog for additional language information and Science-Miner for NER based on Wikidata

```
{"ID": "w109-cb9806ea53",
...
"entities": [
    {
        "label": "citrus fruits",
        "wdid": "Q81513"
    },
...
"language": "https://glottolog.org/resource/languoid/id/kama1365",
"language_iso6393": "woi",
"language_name": "Kamang",
...
"parententities": [
    {
        "label": "Fruit",
        "wdid": "Q1364"
    },
    {
        "label": "Food",
        "wdid": "Q2095"
    }
]
}
```

Model with LIGT





- Word nodes have two labels: RFC 5646 label “en” with a private subtag “-x-lgr” for Leipzig Glossing Rules, following specifications in Section 2.2.7

Present examples: humans

486 results found in 11ms

Page size

5 10 25

Sorting

Relevance

Parent concepts : Animal

Search per field

Vernacular text

Search vernacul

Translation

Search translati

Length (characters)



Length (words)



Filters

Language Iso6393

aey

aqc

beu

[View all](#)

Language name

Amelie

Archi

Bari

[View all](#)

Parent concepts

Food

Organism

Animal

[View more](#)

Concepts

pig

cow

sheep

[View more](#)

Categories

Np

ZCh

a

[View more](#)

Book

A grammar of

Japhug

A grammar of

Mauwake

A grammar of

Rapa Nui

[View more](#)

icq̊a qazō u-kw-nṣye

the.aforementioned sheep 3SG.POSS-SB;PCP-sell

tʰw-kw-ye nu w-pʰe

AOR:DOWNSTREAM-SB;PCP-come[II] DEM 3SG.POSS-DAT
[He told] the person who had come to sell the sheep.' (2003kandZislama) (<https://glottolog.org/resource/languoid/id/japh1234>)

 Language: Japhug

Kum wuel mingrieny tu pelen n-ako.

1SG pig meat PERF dog 3SG.M-eat<3PL>

My pig's meat has been eaten by the dog. ()

 Language: Japhug

Mo ai rō konā hore iho hai 'ārote e pu'a era e
if exist EMPH place cut just_then INS plow IPFV COVER DIST NUM
ono 'o ka va'u rō atu 'uei.
six or CNTG eight EMPH away ox

When a field was ploughed for the first time, it was covered with six or even eight oxen.'
[R539-1.110] (<https://glottolog.org/resource/languoid/id/rapa1244>)

 Language: Rapa Nui

nw tʷ-nu-ndʷm tce tw-mxci smuʷlʷm

DEM IMP-AUTO-take[III] LNK 2-be.rich:FACT prayer

Take (this cattle and, and may you be rich!' (2003kAndzwsqhai2) (<https://glottolog.org/resource/languoid/id/japh1234>)

 Language: Japhug

He haka hāŋai tahī i tū māmoe era.

NTR CAUS feed all ACC DEM sheep DIST

We fed all the sheep.' [R131.008] (<https://glottolog.org/resource/languoid/id/rapa1244>)

 Language: Rapa Nui

› JSON-LD

```
▼ http://purl.org/dc/terms/isPartOf:  
  ▼ 0:  
    @id: https://langsci-press.org/catalog/book/157  
  ▶ http://purl.org/dc/termWors/hasPart:  
  ▶ http://purl.org/dc/terms/license:  
  ▶ http://purl.org/dc/elements/1.1/subject:  
  ▶ https://imtvault.org/content/static/ligt-0.2.ttl#hasWords:  
    ▼ 0:  
      @id: https://imtvault.org/wl09-cb9806ea53_wt"  
      ▶ http://www.w3.org/2000/01/rdf-schema#label:  
      ▼ @type:  
        0: https://purl.org/liodi/ligt#WordTier"  
      ▶ https://imtvault.org/content/static/ligt-0.2.ttl#item:  
        ▼ 0:  
          @id: _:wl09-cb9806ea53_0"  
          ▶ http://www.w3.org/2000/01/rdf-schema#label:  
            ▶ 0:  
            ▶ 1:  
            ▼ @type:  
              ▶ 0: https://imtvault.org/con...static/ligt-0.2.ttl#Word"  
              ▶ https://imtvault.org/content/static/ligt-0.2.ttl#nextWord:  
                ...  
        ▼ 1:  
          @id: _:wl09-cb9806ea53_1"  
          ▶ http://www.w3.org/2000/01/rdf-schema#label:  
            ...
```

present examples: machines

› nquads

```
_:b0 <http://purl.org/dc/terms/identifier> <https://www.wikidata.org/wiki/Q1364>
_:b0 <http://www.w3.org/2000/01/rdf-schema#label> <https://imtvault.org/Fruit> .
_:b1 <http://purl.org/dc/terms/identifier> <https://www.wikidata.org/wiki/Q2095>
_:b1 <http://www.w3.org/2000/01/rdf-schema#label> <https://imtvault.org/Food> .
_:b10 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/woi> .
_:b11 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/en-x-lgr>
_:b12 <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <https://imtvault.org/cont
_:b12 <http://www.w3.org/2000/01/rdf-schema#label> _:b13 .
_:b12 <http://www.w3.org/2000/01/rdf-schema#label> _:b14 .
_:b13 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/woi> .
_:b14 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/en-x-lgr>
_:b2 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/woi> .
_:b3 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/en-x-lgr>
_:b4 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/woi> .
_:b5 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/en-x-lgr> .
_:b6 <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <https://imtvault.org/cont
_:b6 <http://www.w3.org/2000/01/rdf-schema#label> _:b7 .
_:b6 <http://www.w3.org/2000/01/rdf-schema#label> _:b8 .
_:b6 <https://imtvault.org/content/static/ligt-0.2.ttl#nextWord> _:b9 | Nordhoff&Krä
_:b7 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/woi> .
_:b8 <http://purl.org/dc/elements/1.1/language> <https://imtvault.org/en-x-lgr> .
_:b9 <http://www.w3.org/1999/02/22-rdf-syntax-ns#type> <https://imtvault.org/cont
```

Conclusion: added languages



Conclusions

- › There are (semi-)structured sources for linguistic data around
 - › typological databases, endangered language archives, publications
 - › these can be tapped into
 - › expanding the linguistic coverage of the LLOD-cloud
 - › still in the range of 10^2 datapoints, far away from 10 or so required for serious NLP.
- › Glottolog will be a useful to aggregate the data
 - › data is available in CLDF format as well

Thank you

https://imtvault.org/b/85/ex/10-7c7dfc059a.htm 220% Search

IMT VAULT

Example in book 85

Tenk góð wé yu dóñ kán!
thank God SUB 2SG PRF come
Thank God that you have come!'

imtwordsbare [thank, God, SUB, 2SG, PRF, come]
language_glottocode fern1234
language https://glottolog.org/resource/languoid
/id/fern1234
book_ID 85
label Tenk góð wé yu dóñ kán!
wlength 6
license https://creativecommons.org/licenses/by/4.0
book_metalanguage eng
language_iso6393 fpe
book_title A grammar of Pichi
book_URL https://langsci-press.org/catalog/book/85
ID 10-7c7dfc059a
categories [prf, sub]
language_name Pichi
srcwordsbare [Tenk, góð, wé, yu, dóñ, kán!]
clength 57