

# Wortschatz / Leipzig Corpora Collection

## Working with under-resourced languages

# Collecting smaller languages

## Problems

- Often, only a small part of the Web in country is in the language under consideration.
- More resources can be found in .com domains, but they are difficult to identify because of the size of the .com TLD.

## Criteria

- More than 1.000.000 speakers
- Less than 1.000.000 sentences found so far

- Wordlist based approach using wordlists from
  - Universal Declaration of Human Rights (UDHR)
  - Watchtower documents
- Tuples of words for search engine queries
- Usage of Webservice-based APIs – if available - (Google, Microsoft Bing) to collect URLs
  - Can be iterated

Language	Code	Country	Speakers	Sentences
Somali	<a href="#">som</a>	Somalia	13871700	121545
Swahili	<a href="#">swa</a>	Tanzania	15458390	111927
Zulu	<a href="#">zul</a>	South Africa	10349100	75869
<a href="#">Shona</a>	<a href="#">sna</a>	Zimbabwe	10759200	52060
Xhosa	<a href="#">xho</a>	South Africa	7817300	41983
Amharic	<a href="#">amh</a>	<a href="#">Ethiopia</a>	17528500	24856
Northern Sotho	<a href="#">nso</a>	South Africa	4101000	19000
Southern Sotho	<a href="#">sot</a>	South Africa, Lesotho	6024000	18000
Yoruba	<a href="#">yor</a>	Nigeria	19380800	12886
Ganda	<a href="#">lug</a>	Uganda	4130000	11058
Southern Ndebele	<a href="#">nbl</a>	South Africa	1090000	8000
Wolof	<a href="#">wol</a>	Senegal	3976500	5715
Hausa	<a href="#">hau</a>	Nigeria	24988000	5506
Tswana	<a href="#">tsn</a>	South Africa, Botswana	4521700	5229
Tsonga	<a href="#">tso</a>	South Africa	3669000	3264
<a href="#">Lingala</a>	<a href="#">lin</a>	Dem. Rep. of Congo	2141300	2783
<a href="#">Ndonga</a>	<a href="#">ndo</a>	Namibia	1070000	2659

<b>Language</b>	<b>Code</b>	<b>Country</b>	<b>Speakers</b>	<b>Sentences</b>
Venda	<u>ven</u>	South Africa	1064000	1893
<u>Igbo</u>	<u>ibo</u>	Nigeria	18000000	1716
<u>Bamanankan</u>	<u>bam</u>	Mali	2772340	1643
Oromo	<u>orm</u>	<u>Ethiopia</u>	17467900	1442
<u>Éwé</u>	ewe	Ghana	3112000	1000
Swati	<u>ssw</u>	South Africa, Swaziland	2034200	901
<u>Rundi</u>	run	Burundi	4851000	826
Tigrigna	tir	Ethiopia, Eritrea	5791710	805
<u>Fulah</u>	<u>ful</u>	Cameroon	22246400	519
Akan	aka	Ghana	8314600	505
<u>Koongo</u>	<u>kng</u>	Dem. Rep. of Congo	5000000	490
Rwanda	<u>kin</u>	Rwanda	7504900	469
Nyanja	<u>nya</u>	Malawi	8659700	446
<u>Tumbuka</u>	<u>tum</u>	Malawi	1142000	300
<u>Acholi</u>	<u>ach</u>	Uganda	1215000	192
<u>Gikuyu</u>	<u>kik</u>	<u>Kenya</u>	6623000	138
<u>Pulaar</u>	<u>fuc</u>	Senegal	3691000	18

Language	Code	Country	Speakers
<a href="#">Umbundu</a>	<a href="#">umb</a>	<a href="#">Angola</a>	4002880
<a href="#">Kimbundu</a>	<a href="#">kmb</a>	<a href="#">Angola</a>	4000000
<a href="#">Fon</a>	<a href="#">fon</a>	<a href="#">Benin</a>	1435500
<a href="#">Mòoré</a>	<a href="#">mos</a>	<a href="#">Burkina Faso</a>	5061700
<a href="#">Jula</a>	<a href="#">dyu</a>	<a href="#">Burkina Faso</a>	1229000
<a href="#">Dan</a>	<a href="#">dnj</a>	<a href="#">Côte d'Ivoire</a>	1610800
<a href="#">Baoulé</a>	<a href="#">bci</a>	<a href="#">Côte d'Ivoire</a>	2130000
<a href="#">Zande</a>	<a href="#">zne</a>	<a href="#">Dem. Rep. of Congo</a>	1142000
<a href="#">Yombe</a>	<a href="#">yom</a>	<a href="#">Dem. Rep. of Congo</a>	1056400
<a href="#">Songe</a>	<a href="#">sop</a>	<a href="#">Dem. Rep. of Congo</a>	1000000
<a href="#">Ngbaka</a>	<a href="#">nga</a>	<a href="#">Dem. Rep. of Congo</a>	1016650
<a href="#">Luba-Katanga</a>	<a href="#">lub</a>	<a href="#">Dem. Rep. of Congo</a>	1510000
<a href="#">Luba-Kasai</a>	<a href="#">lua</a>	<a href="#">Dem. Rep. of Congo</a>	6300000
<a href="#">Kongo</a>	<a href="#">kon</a>	<a href="#">Dem. Rep. of Congo</a>	5644100
<a href="#">Kituba</a>	<a href="#">ktu</a>	<a href="#">Dem. Rep. of Congo</a>	4200000
<a href="#">Chokwe</a>	<a href="#">cjk</a>	<a href="#">Dem. Rep. of Congo</a>	1009780
<a href="#">Alur</a>	<a href="#">alz</a>	<a href="#">Dem. Rep. of Congo</a>	1367000
<a href="#">Tigré</a>	<a href="#">tig</a>	<a href="#">Eritrea</a>	1050000
<a href="#">Wolaytta</a>	<a href="#">wal</a>	<a href="#">Ethiopia</a>	1710000
<a href="#">Sidamo</a>	<a href="#">sid</a>	<a href="#">Ethiopia</a>	2980000
<a href="#">Hadiyya</a>	<a href="#">hdy</a>	<a href="#">Ethiopia</a>	1250000
<a href="#">Gamo</a>	<a href="#">gmv</a>	<a href="#">Ethiopia</a>	1110000
<a href="#">Afar</a>	<a href="#">aar</a>	<a href="#">Ethiopia</a>	1078200
<a href="#">Abron</a>	<a href="#">abr</a>	<a href="#">Ghana</a>	1182000
<a href="#">Susu</a>	<a href="#">sus</a>	<a href="#">Guinea</a>	1060280
<a href="#">Pular</a>	<a href="#">fuf</a>	<a href="#">Guinea</a>	2929200
<a href="#">Kpelle</a>	<a href="#">kpe</a>	<a href="#">Guinea</a>	1220000
<a href="#">Fang</a>	<a href="#">fan</a>	<a href="#">Guinea</a>	1027900
<a href="#">Eastern Maninkakan</a>	<a href="#">emk</a>	<a href="#">Guinea</a>	3531800
<a href="#">Oluluyia</a>	<a href="#">luy</a>	<a href="#">Kenya</a>	5199727
<a href="#">Maasai</a>	<a href="#">mas</a>	<a href="#">Kenya</a>	1045000
<a href="#">Lubukusu</a>	<a href="#">bvk</a>	<a href="#">Kenya</a>	1433000
<a href="#">Kipsigis</a>	<a href="#">sgc</a>	<a href="#">Kenya</a>	1916000
<a href="#">Kimĩiru</a>	<a href="#">mer</a>	<a href="#">Kenya</a>	1658000
<a href="#">Kamba</a>	<a href="#">kam</a>	<a href="#">Kenya</a>	3893000
<a href="#">Kalenjin</a>	<a href="#">kln</a>	<a href="#">Kenya</a>	5123400
<a href="#">Ekegusii</a>	<a href="#">guz</a>	<a href="#">Kenya</a>	2120300
<a href="#">Dholuo</a>	<a href="#">luo</a>	<a href="#">Kenya</a>	4410000
<a href="#">Yao</a>	<a href="#">yao</a>	<a href="#">Malawi</a>	1916000
<a href="#">Soninke</a>	<a href="#">snk</a>	<a href="#">Mali</a>	1250000
<a href="#">Maasina Fulfulde</a>	<a href="#">ffm</a>	<a href="#">Mali</a>	1008500
<a href="#">Tswa</a>	<a href="#">tsc</a>	<a href="#">Mozambique</a>	1180000

Language	Code	Country	Speakers
<a href="#">Sena</a>	<a href="#">seh</a>	<a href="#">Mozambique</a>	1340000
<a href="#">Makhuwa-Meetto</a>	<a href="#">mgh</a>	<a href="#">Mozambique</a>	1348000
<a href="#">Makhuwa</a>	<a href="#">vmw</a>	<a href="#">Mozambique</a>	3090000
<a href="#">Lomwe</a>	<a href="#">ngl</a>	<a href="#">Mozambique</a>	1500000
<a href="#">Zarma</a>	<a href="#">dje</a>	<a href="#">Niger</a>	2438400
<a href="#">Tamashek</a>	<a href="#">tmh</a>	<a href="#">Niger</a>	1248200
<a href="#">Tiv</a>	<a href="#">tiv</a>	<a href="#">Nigeria</a>	2210000
<a href="#">Nigerian Pidgin</a>	<a href="#">pcm</a>	<a href="#">Nigeria</a>	30000000
<a href="#">Nigerian Fulfulde</a>	<a href="#">fuv</a>	<a href="#">Nigeria</a>	11500000
<a href="#">Kanuri</a>	<a href="#">kau</a>	<a href="#">Nigeria</a>	3760500
<a href="#">Izon</a>	<a href="#">ijc</a>	<a href="#">Nigeria</a>	1000000
<a href="#">Ibibio</a>	<a href="#">ibb</a>	<a href="#">Nigeria</a>	1750000
<a href="#">Edo</a>	<a href="#">bin</a>	<a href="#">Nigeria</a>	1000000
<a href="#">Ebira</a>	<a href="#">igb</a>	<a href="#">Nigeria</a>	1000000
<a href="#">Central Kanuri</a>	<a href="#">knc</a>	<a href="#">Nigeria</a>	3240500
<a href="#">Berom</a>	<a href="#">bom</a>	<a href="#">Nigeria</a>	1000000
<a href="#">Anaang</a>	<a href="#">anw</a>	<a href="#">Nigeria</a>	1400000
<a href="#">Serer-Sine</a>	<a href="#">srr</a>	<a href="#">Senegal</a>	1161900
<a href="#">Mandinka</a>	<a href="#">mnk</a>	<a href="#">Senegal</a>	1346000
<a href="#">Mandingo</a>	<a href="#">man</a>	<a href="#">Senegal</a>	6496300
<a href="#">Themne</a>	<a href="#">tem</a>	<a href="#">Sierra Leone</a>	1230000
<a href="#">Mende</a>	<a href="#">men</a>	<a href="#">Sierra Leone</a>	1499700
<a href="#">Dinka</a>	<a href="#">din</a>	<a href="#">Sudan</a>	1365900
<a href="#">Bedawiyet</a>	<a href="#">bej</a>	<a href="#">Sudan</a>	1186000
<a href="#">Sukuma</a>	<a href="#">suk</a>	<a href="#">Tanzania</a>	5430000
<a href="#">Nyakyusa-Ngonde</a>	<a href="#">nyy</a>	<a href="#">Tanzania</a>	1105000
<a href="#">Makonde</a>	<a href="#">kde</a>	<a href="#">Tanzania</a>	1340000
<a href="#">Haya</a>	<a href="#">hay</a>	<a href="#">Tanzania</a>	1300000
<a href="#">Gogo</a>	<a href="#">gog</a>	<a href="#">Tanzania</a>	1440000
<a href="#">Teso</a>	<a href="#">teo</a>	<a href="#">Uganda</a>	1849000
<a href="#">Soga</a>	<a href="#">xog</a>	<a href="#">Uganda</a>	2060000
<a href="#">Nyankore</a>	<a href="#">nyn</a>	<a href="#">Uganda</a>	2330000
<a href="#">Masaaba</a>	<a href="#">myx</a>	<a href="#">Uganda</a>	1120000
<a href="#">Lugbara</a>	<a href="#">lgb</a>	<a href="#">Uganda</a>	1637000
<a href="#">Lango</a>	<a href="#">laj</a>	<a href="#">Uganda</a>	1490000
<a href="#">Chiga</a>	<a href="#">cgg</a>	<a href="#">Uganda</a>	1580000
<a href="#">Bemba</a>	<a href="#">bem</a>	<a href="#">Zambia</a>	3602000
<a href="#">Tonga</a>	<a href="#">toj</a>	<a href="#">Zambia, Zimbabwe</a>	1127000
<a href="#">Ndebele</a>	<a href="#">nde</a>	<a href="#">Zimbabwe</a>	1572800
<a href="#">Ndau</a>	<a href="#">ndc</a>	<a href="#">Zimbabwe</a>	2380000
<a href="#">Manyika</a>	<a href="#">mxc</a>	<a href="#">Zimbabwe</a>	1025000

- Portal to allow external support
- User friendly interface, but work in progress
- Try it: <http://curl.corpora.uni-leipzig.de>

Collecting Web Pages for

# Under-Resourced Languages

On this website you can contribute to corpus collection for under-resourced languages by simply entering a URL. The URLs or domains you provide will be crawled and reviewed for text data in the respective language. After processing you will be presented with statistics for the URLs you provided. The corpora created will be free for download.

Thank you for your Support!



Step 1 » Please select a language.

- Yombe [yom] - Dem. Rep. of Congo
- Yoruba [yor] - Nigeria
- Zande [zne] - Dem. Rep. of Congo
- Zarma [dje] - Niger
- Zaza [zza] - Turkey
- Zhuang [zha] - China
- Zulu [zul] - South Africa

Continue





Processing Data



## Step 2 » Please insert your URLs

Please add one URL per line!

Upload a File!

Previous

Continue

## Job details for amh-2016-09-20-07-51-19

**Job Name** amh-2016-09-20-07-51-19  
**Submitted on** 2016-09-20 07:51:20  
**Processed** 2016-09-20 14:04:04 / 2016-09-20 20:20:11  
**Language** [Amharic \[amh\]](#)  
**# URLs** 8  
**Status** All the URLs have been processed.  
No details available.

URL List					
Show	10	▼	entries		
			Search: <input type="text"/>		
Status	⌵	Language	⌴	URL	⌴
<input checked="" type="checkbox"/> finished		amh		<a href="http://www.ebc.et/web/news">http://www.ebc.et/web/news</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="http://ethiopiazare.com/">http://ethiopiazare.com/</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="https://www.lds.org/general-conference/201...">https://www.lds.org/general-conference/201 ...</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="http://www.formerethiopianairforce.com/">http://www.formerethiopianairforce.com/</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="http://etntc.org/">http://etntc.org/</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="http://amharic.voanews.com/a/ethiopians-de...">http://amharic.voanews.com/a/ethiopians-de ...</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="http://www.ginbot7.org/">http://www.ginbot7.org/</a>	
<input checked="" type="checkbox"/> finished		amh		<a href="http://www.ethioshengo.org/">http://www.ethioshengo.org/</a>	

# Corpus Information for Amharic [amh] Ethiopia

**Language** Amharic

**ISO Code** [amh](#) 

**Country** Ethiopia

**Corpus Name** amh\_community\_2016

**Tokens** 637583

**Types** 133454

**Sentences** 43083

**Sources (URLs)** 7662

**Build date** 2016-09-20

**URLs** [List of URLs](#) download

[List of Domains](#) download

**Download** [amh\\_community\\_2016](#) 2016-09-20

**Contact** No contact person for this language.

Use this [Contact](#)  to add contact details.



# CORPORA COLLECTION

UNIVERSITY LEIPZIG

[Home](#)[Submissions / Process Status](#)[Corpora Downloads](#)<http://am.wikipedia.org/>

33672

<http://www.ginbot7.org/>

2212

<http://etntc.org/>

1738

<http://www.ethioshengo.org/>

1392

<http://www.ebc.et/>

1063

<https://www.lds.org/>

972

<http://www.formerethiopianairforce.com/>

443

<http://amharic.voanews.com/>

104

Statistics for user hau-2016-05-23-07-43-24  
-----

We used Heritrix to crawl your URLs.

You submitted 1 seed URLs. 1 seeds could be crawled but 0 couldn't for various reasons.

Heritrix downloaded 6814 URLs from which 6589 URLs were valid.

In total we downloaded 129 MiB in 6h4s10ms. We had a download speed of 0.3 URI/sec.

We then extracted the text content from the webpages. Following this, we separated 3144 documents by language using linguistic methods like trigrams and stop words, and the encoding.

For your language 'hau' we found 1015 documents, 32.28 % of the total. Furthermore, 6 other languages were found but discarded.

After the preprocessing, the documents were split in sentences and the sentences classified by language.

For the language 'hau' we found 8800 sentences from a total of 9057.

We successfully re-created the corpus 'hau\_community\_2016' with your additional language data. The URLs you provided added 805 source URLs, 8431 sentences and 4845 words. The old corpus had 1687 sources, 33158 sentences and 35387 words.

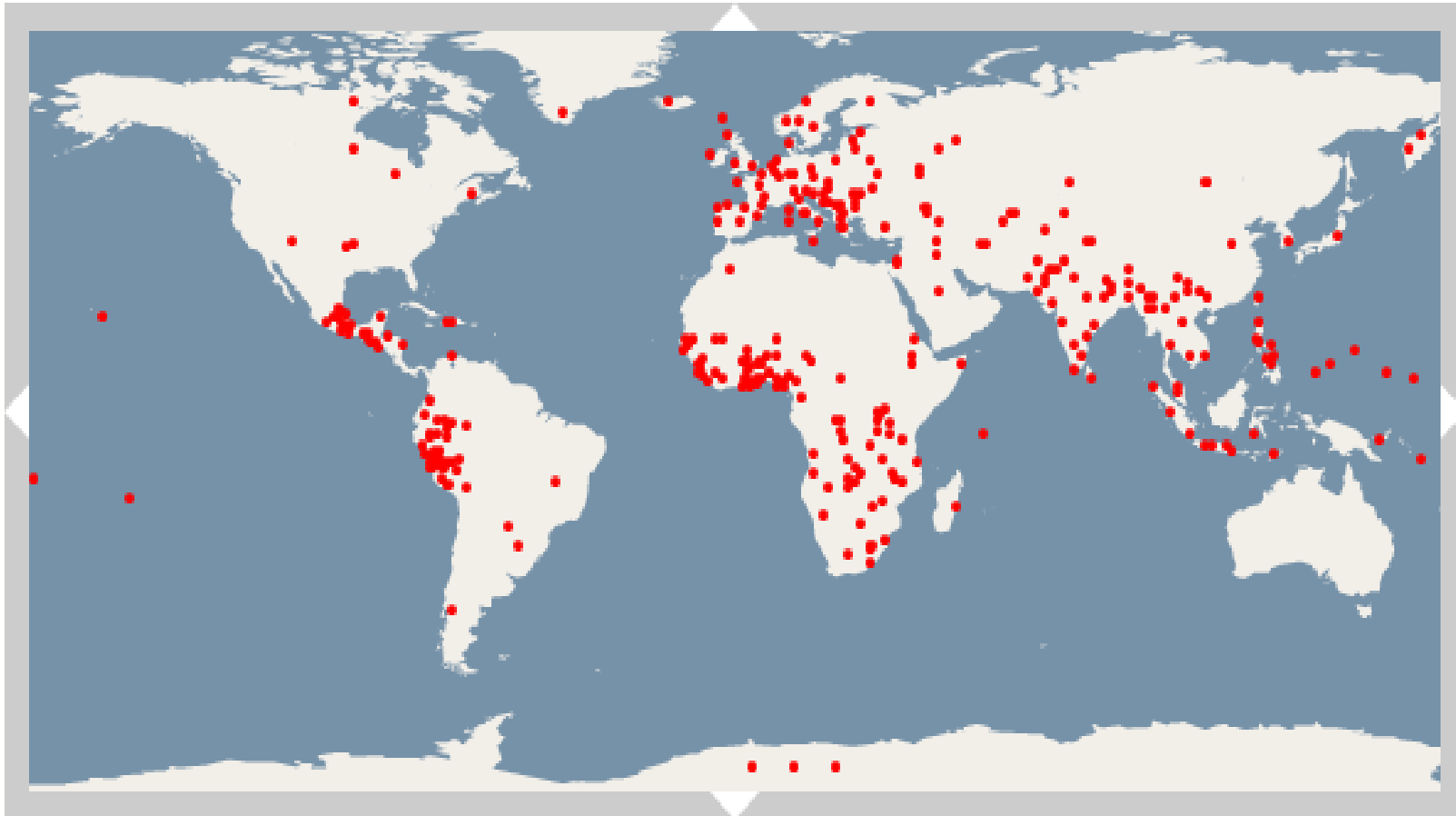
---

# Live Demo

In two steps:

1. step: Identification of the document language using the following algorithm:  
Use lists of the 50 most frequent words for many languages. The language having by far the most hits wins. Documents are stored according to their language.
2. step: The same for each sentence using 5.000 words per sentence. Sentences not in the expected language are rejected.

- The **Universal Declaration of Human Rights** is available in more than 360 languages in UNICODE
- Size of the German version: approx. 1700 words

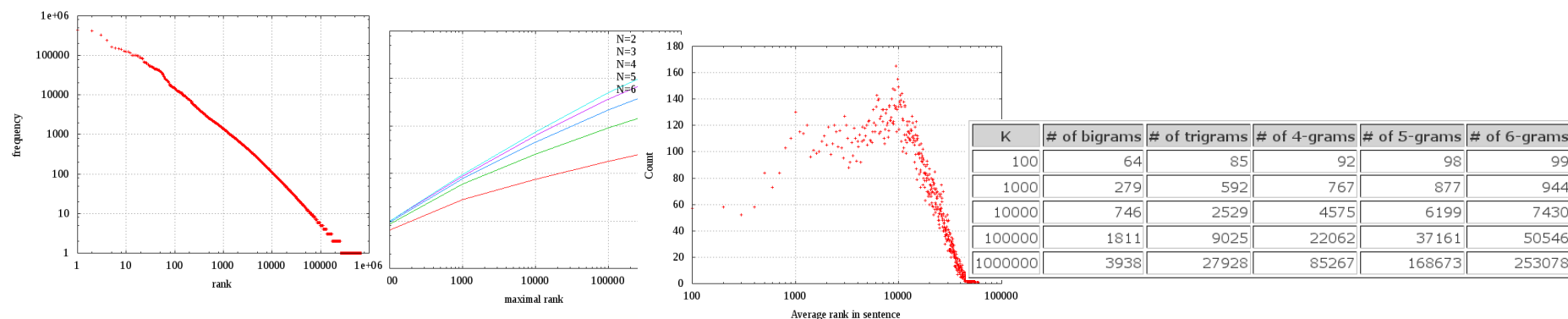




## Bibles (or parts of the Bible)

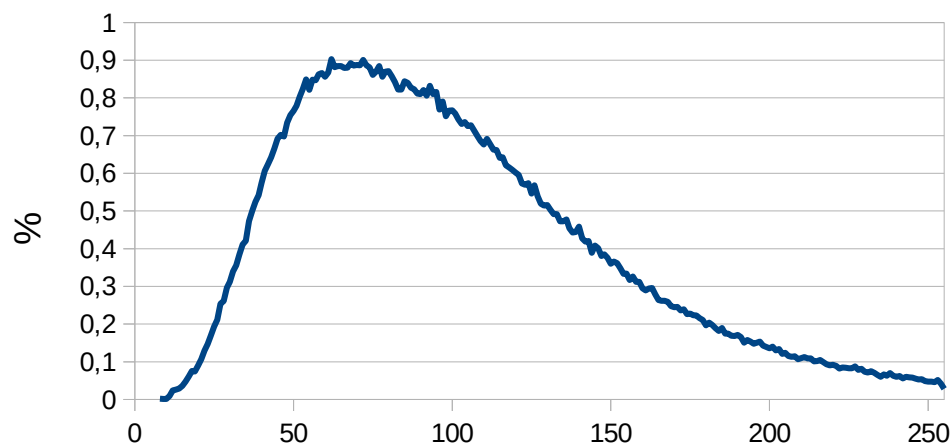
- in nearly 1000 languages
- With typically around 7000 verses
- Sources: [bible.is](http://bible.is) and other websites

- Enrichment of corpora with statistical annotations
  - Word frequencies
  - Co-occurrence frequencies and significance (based on left or right neighbours or whole sentences)
- Creation of further corpora statistics
  - 230 statistical properties based on letter, word, sentence, sources level etc.

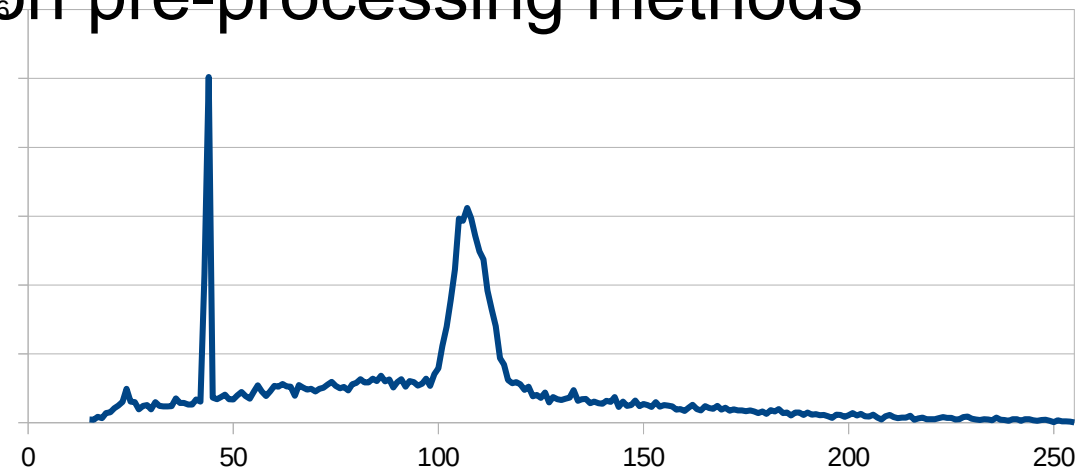


# Applications of Corpora Statistics

- Automatic analysis of statistical properties for quality enhancements
  - Usage of typical distributions of linguistic properties
    - Length distributions (word, sentence, paragraph)
    - Character and n-gram distributions etc.
  - Basis for further cleaning procedures
  - Basis for adaptations on pre-processing methods



Hindi newspaper 2010



Sundanese Wikipedia 2007

- Correlations between parameters

corpus-based  
properties  
(measured)

typological  
parameters

- Determine correlations:

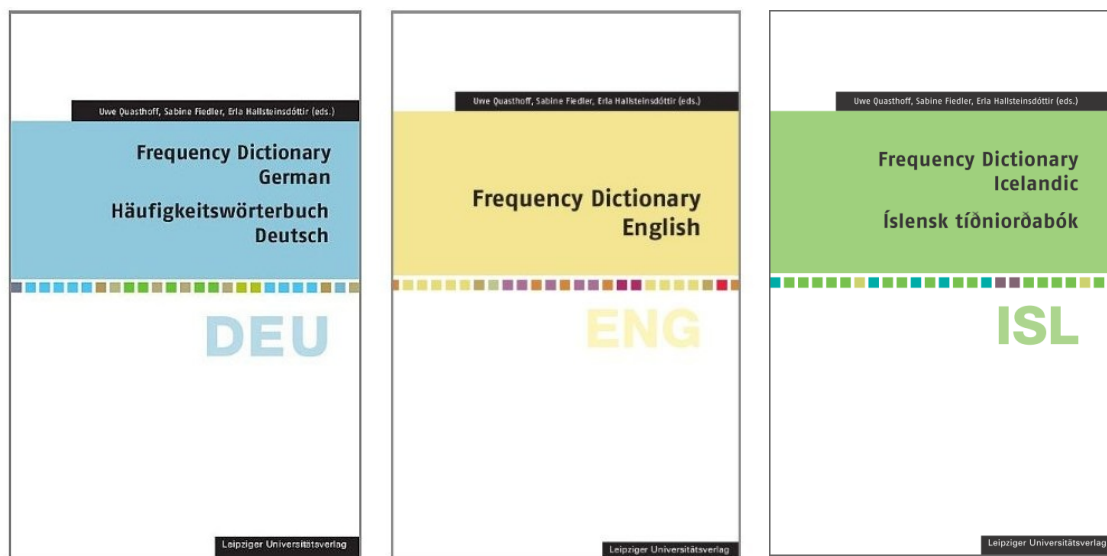
- Between different measurements
- Between measured properties and typological parameters
- Typological parameters:
  - Refer e.g. to morphology or syntax
  - Source: World Atlas of Language Structures (WALS)
- Usage of quantitative Methods:
  - Correlation Analysis
  - Tests of Significance

# Comparison of Language Parameters

Lang	Alphabet	VC Ratio	Cap in text	Cap in top100K	Zipf slope	Cov: 10	Cov: 100	Cov: 1000	Len in text	Len in top100K
<b>CES</b>	<b>42</b>	<b>42:58</b>	<b>13.9%</b>	<b>24,3%</b>	<b>-0.96</b>	<b>13.5%</b>	<b>28.6%</b>	<b>47.0%</b>	<b>5.37</b>	<b>9.08</b>
<b>DEU</b>	31	38:62	36.0 %	68,3%	-1.07	15.5 %	39.1 %	59.2 %	5.86	9.73
<b>ENG</b>	26	40:60	17,6%	48,6 %	-1.10	21.4%	41.5%	63.2%	4.96	7.80
<b>FRA</b>	ca. 41	44:56	12.9%	40,6 %	-1.12	22.1%	46.8%	66.1%	4.82	7.80
<b>HUN</b>	ca. 38	42:58	12.9%	16,9 %	-0.95	19,7%	31,0%	47,9%	6,27	8.96
<b>IND</b>	26	42:58	19.5%	49.9 %	-1.08	13.5%	31.3%	58.8%	6.08	9.28
<b>ISL</b>	ca. 34	39:61	13,5%	30,0 %	-1.06	21,5%	40.3%	60.5%	5.08	8.47
<b>UKR</b>	34	46:54	16.3%	38.3%	-1.14	23.1%	46.1%	69.1%	4.36	8.48

- Correlations between measurements:
  - Negative Correlation between length of words in characters and sentence length in words (Cor=-0,55,  $p < 0,001\%$ )
  - The more syllables per word the fewer words per sentence (Cor=-0,49,  $p < 0,001\%$ )
  - The more syllables per word, the more syllables per sentence (Cor=0,47,  $p < 0,001\%$ )
- Correlations between measurements and typological parameters:
  - Morphological type: concatenative vs. isolating
    - Isolating languages have shorter words:  $p < 1\%$ , average of 8.43 and 6.95
  - Word order: SOV vs. SVO
    - SOV languages typically use more syllables per word:  $p < 0,005\%$ , average of 2.21 and 1.89

- Book series, published by Leipzig University Press since 2011
- Most frequent words of a language with frequency information
  - Printed: Top 10,000 words
  - Electronically (CD): Up to top 1,000,000 words



- Searching for language experts for further languages
- [quasthoff@informatik.uni-leipzig.de](mailto:quasthoff@informatik.uni-leipzig.de)



Additional statistical information on word list about:

- alphabet, distribution of vowels and consonants
  - word length (average and for different frequency ranges)
  - Number of syllables per word
  - Longest words for different frequency ranges
  - Corpus timeline
-

## In preparation:

- **NLD** – Dutch
- **VIE** – Vietnamese
- **CES** – Czech
- **POL** – Polish
- **RUS** – Russian
- **GLG** – Galician
- **DAN** – Danish
- **SWE** – Swedish
- **AFR** – Afrikaans
- **ZUL** – Zulu

