

# Language Comparison with CLDF

Robert Forkel

Department of Linguistic and Cultural Evolution  
MPI EVA, Leipzig

# Language Comparison in “Analog” Times

# The "Pallas dictionary" from 1787/89

I. Б О Г Ъ.		
1 — По Славянски	-	Богъ.
2 — Славяно - Венгер-	ски	Бугъ.
3 — Иалиръски	-	Боѓъ.
4 — Богемски	-	Бу.
5 — Сербски	-	Богъ.
6 — Вендски	-	Богъ.
7 — Сорабски	-	Богъ.
8 — Подабски	-	Бусацъ.
9 — Кашубски	-	Богъ.
10 — Польски	-	Богъ.
11 — Малороссийски	-	Бигъ.
12 — Сузdalъски	-	Стодъ.
13 — Кельшки	-	Діу, Ю.
14 — Бремански	-	Дуз, Доз.
15 — Басконски	-	Дувъ, Юнъ, Янне, Яинкоа, Энкоа.
16 — Ирландски	-	Ия.
17 — Эрзо - Шотланд-	ски	Дїа.
18 — Валски	-	Дїу.
19 — Корнвалски	-	Деу.
20 — Елински	-	Ѳеосъ.
21 — Ново-Гречески	-	Ѳеосъ.
22 — Лашински	-	Деусъ.
23 — Ишалански	-	Діо.
24 — Неаполитански	-	Дадо.
25 — Испански	-	Діосъ.
26 — Португальски	-	Деосъ.
27 — Роменски и дре-		вне-Французски
		Деу, Декъ, Ді- едъ, Діоръ.
34 — Понижне-Германски	Годъ.	
35 — Германски	-	Готицъ.
36 — Цимбрски	-	Гиппицъ.
37 — Датски	-	Гудъ.
38 — Исландски	-	Гудъ.
39 — Шведски	-	Гудъ.
40 — Голландски	-	Годъ.
41 — Фризски	-	Годъ.
42 — Литовски	-	Дізвасъ.
43 — Латышски	-	Даєсъ.
44 — Кривинго - Ли-		вонски
		Дізвасъ.
45 — Албански	-	Перенди.
46 — Волошки	-	Думнезеу.
47 — Венгерски	-	Иштень.
48 — Аварски	-	Бечасъ.
49 — Кубачински	-	Бечасъ.
50 — Лезгински, рода		
		Анцуғъ
		Бедшетъ.
51 — — р. Джаръ	-	Бедшетъ.
52 — — р. Хунзагъ	-	Беджетъ.
53 — — р. Дидо	-	Бедшешъ.
54 — Чюхонски	-	Юмала.
55 — Эстландски	-	Юммалъ.
56 — Корельски	-	Юмала.
57 — Олонецки	-	Юмаль.
58 — Лопарски	-	Юбмелъ, Ибмелъ
59 — Зырянски	-	Іенъ.
60 — Пермакски	-	Іенъ-хэнъ.
61 — Мордовски	-	Паасъ.
62 — Мокшански	-	Шкай, Шкипаасъ
63 — Черемиски	-	Юму.

## I. БОГЪ.

- |                    |           |                             |
|--------------------|-----------|-----------------------------|
| 1 По Славянски -   | Богъ.     | 34 По Нижне-Германски Годъ. |
| 2 —Славяно-Венгер- |           | 35 —Германски - Готий.      |
|                    | ски Бугъ. | 36 —Цимбрски - Гипишъ.      |
| 3 —Иллирйски       | Боогъ.    | 37 —Датски - Гудъ.          |
| 4 —Богемски        | Бу.       | 38 —Исландски - Гудъ.       |

88

## God (60).

Number in General List.		Number in General List.		Number in General List.	
<b>AGGLUTINATIVE NON-INDIAN LANGUAGES.</b>					
Japanese . . .	<i>Kami</i>	35. Pwo, literary . .	<i>-Ka -sā /yuwa</i>	123. Abor . . . . .	
Ainu . . .	<i>Kamui</i>	" Bassein . .	<i>Yuə, Chai'</i>	124. Miri . . . . .	
Korean . . .	<i>Hanñim</i>	" Maulmein . .	<i>Yuə, Cha'</i>	125. Dafla . . . . .	
Turki . . .	<i>Tiŋri, Auŋan, Bir-i-bär</i>	36. Taungθu . .	<i>P'ra</i>	126. Mišmi, Digārū . .	
Manchu . . .	<i>Abkă-i ejen</i> (Lord of heaven)	34. Sgă, literary . .	<i>-Ka uə</i>	<i>Nin-ya (?)</i>	
Mongolian . . .	<i>Tegri</i>	" spoken . .	<i>G'să yuă</i>	<i>Miju . . . Se-lap</i>	
Saukpă . . .	...	32. Bwè . .	<i>Bi je uă</i>	<i>Lolo-Mos'o Group.</i>	
Basque . . .	<i>Jaungoiko, Jainko, Dzipo</i>	41a. Wewaw . .	<i>G'să ȣiək</i>	Si-hia . . . <i>K-num, q-num</i> (Heaven)	
		" <i>Wewankwe</i>	<i>ȣiək</i>	273. Lolo, /N <sup>i</sup> . . . <i>/M" uə /p'a</i>	
				A-hi . . . <i>Mu<sup>o</sup> sa<sup>o</sup> p'o<sup>o</sup></i>	

# Data curation requires "comparability", too.

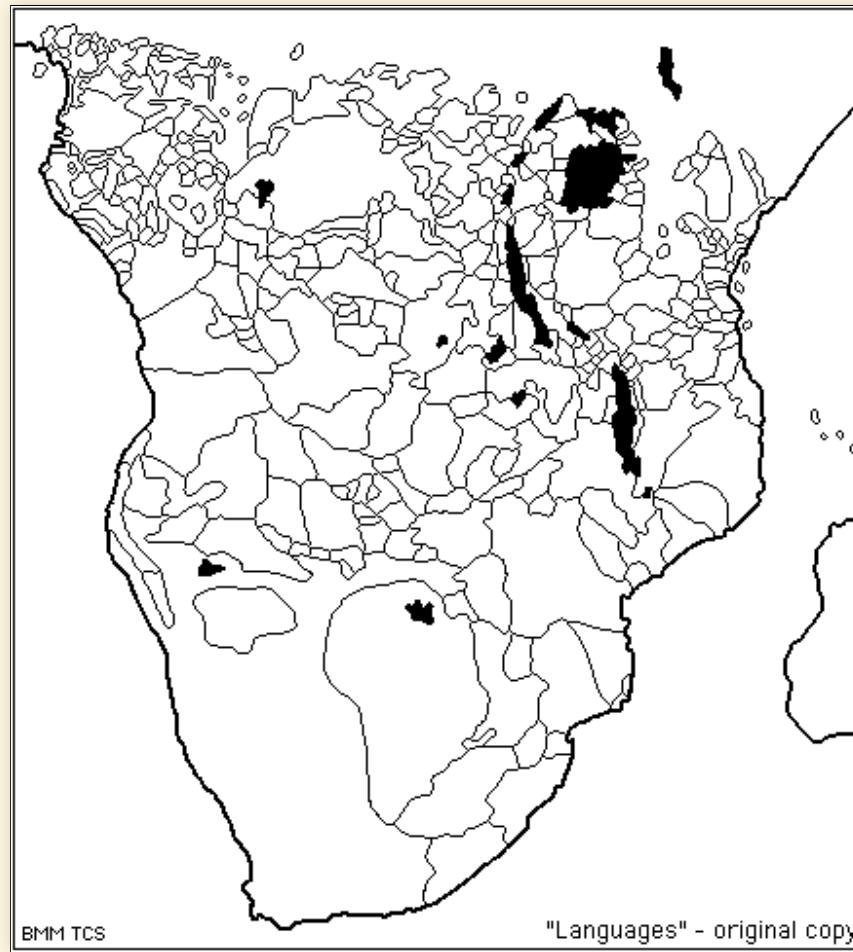
201. морюк-	244. урана - мусунчуканын
202. чирчи-	245. януу -
203. смычко - яасанын	246. яаги -
204. юрга-	247. кыпч - кысатычын
205. морено -	248. пылг - нааманын түшрекин
206. тедоганыз - юнсан	249. моне - яадын
207. сонгы - юнсан	250. жапа - мынчилек
208. мороз -	251. зама -
209. супор -	252. оюнчан -
210. муро -	253. дээдэ - яадынхана
211. ногы - яанык	254. ноги -
212. наул -	255. нога -
213. ына -	256. нога - енишес
214. ногы -	257. нога - одыма
215. зано -	258. нога - алжарас
216. шылб -	259. ноги -
217. энди -	260. ноги -
218. ынба - <sup>томи</sup> <sub>жыныс</sub>	261. наал - яадынхана
219. наал - тендерин	262. ноги -
220. күрдэ -	263. ноги -
221. 220.чаны - <sup>жарын</sup> <sub>жанынанын</sub>	264. то -
222. аисома -	265. тоот -
223. шурмана - тарытесеми	266. тоопт -
224. матары -	267. тоот -
225. дына - мүшнинчуми.	268. тоо -
226. балыкчи - наадызас	269. тоо -
227. дынад -	270. тоо -
228. синер -	271. тоо -
229. 24/16 - <sup>ядалык</sup> <sub>Дончакыра</sub>	272. тоо -
230. оотина - тарадычынне	273. тоо -
231. зорын -	274. тоо -
232. кичи -	275. тоо -
233. кичиц -	276. тоо -
234. кичи - <sup>кызы</sup> <sub>кызы</sub>	277. тоо -
235. 26/18 -	278. тоо -
236. дынан - дысалдын	279. тоо -
237. оотын - немүнчуми	280. тоо - яончы калориянык
238. оотын - яанык	281. тоо -

# Going “Digital”

Nande.Kavutirwaki1978.txt  
~/Downloads

1	Prefix	Stem	POS	Class	Gloss	Word	Tone
2		-abuü	n.	7/8	le fumier, le d <small>ə</small> potoir; g <small>ə</small> n <small>ə</small> ralement derri <small>ə</small> re la hutte et o <small>ə</small> 9		
	l'on jette les ordures				eky <small>ə</small> buü	A	
3		-abuü	n.	14/6	la bi <small>ə</small> re; toute boisson alcoolis <small>ə</small> e import <small>ə</small> e ou fabriqu <small>ə</small> e	base	
	de produits locaux				bwabuü	D	
4		-agaliü-		v.	trans. faire souffrir	eryagaly <small>ə</small> 9	A
5		-agalo-	v.	intr.	souffrir	la suite d'une preuve physique ou morale	
	eryagalw <small>ə</small> 89				A		
6		-ag <small>ə</small> nd <small>ə</small> 87	n.	7/8	la case de r <small>ə</small> union pour les anciens du village		8E
	kyag <small>ə</small> nd <small>ə</small> 87				D		
7		-ahul-	v.	trans.	nommer; dire le nom de quelqu'un ou de quelque chose		
	ery <small>ə</small> h <small>ə</small> la				B		
8		-ak-	v.	intr.	tre allum <small>ə</small> ; tre br <small>ə</small> l <small>ə</small>	eryak <small>ə</small> 87	A
9		-aka	n.	3/4	ensemble de plantes se trouvant dans un champ		
10		-akakala	n.	7/8	la lumi <small>ə</small> re; un faisceau de lumi <small>ə</small> re		
					D	mwaka	D
11		-ala	n.	8	la main	by <small>ə</small> la	B
12		-ala	n.	7/8	l'ongle des doigts ou des orteils		
						ky <small>ə</small> la	B
13		-li	n.	1/2	la fille de	omw <small>ə</small> li L-	
14		-ly	n.	8	la nourriture	eby <small>ə</small> ly <small>ə</small>	L-
15		-alyana	n.	1/2	la belle-soeur (tant du c <small>ə</small> t <small>ə</small> de l'homme que de la		
	femme)					D	
16		-amba	n.	7/8	sorte de fruit comestible se trouvant dans la for <small>ə</small> t		8E

Plain Text ▾ Tab Width: 8 ▾ Ln 1, Col 1 ▾ INS



title	Damal
fill	#51CC70
family	null
cldf:languageRef	dama1272
fill-opacity	0.8

+ Add row    Add simplestyle properties

Properties    Info

Save    Cancel    Delete feature

```

1 {  

2   "type": "FeatureCollection",  

3   "features": [  

4     {  

5       "type": "Feature",  

6       "properties": {  

7         "title": "Hlai",  

8         "fill": "#51CC70",  

9         "family": "Tai-Kadai",  

10        "cldf:languageReference":  

11          "fill-opacity": 0.8  

12      },  

13      "geometry": {  

14        "type": "Polygon",  

15        "coordinates": [ ]  

16      }  

17    },  

18    {  

19      "type": "Feature",  

20      "properties": {  

21        "title": "Tsat",  

22        "fill": "#CC5151",  

23        "family": "Austronesian",  

24        "cldf:languageReference":  

25          "fill-opacity": 0.8  

26      },  

27      "geometry": {  

28        "type": "Polygon",  

29        "coordinates": [ ]  

30      }  

31    }  

32  }  

33 }

```

# Going On-Line

← → ⌂ https://pallas.ivdnt.org/lexit2/?db=pallas&lang=en ⭐

/instituut voor de Nederlandse taal/ The Digital Pallas Peter Simon Pallas, *Comparative Dictionary of All Languages and Dialects* (1790-1791)

Pallas:  
61.970 row(s) found  
Show 10 rows

Search whole table:

First Previous 1 2 3 4 5 ... 6197 Next Last

word	modern word	transliteration	Pallas concept	modern concept	English concept	Pallas language	modern language	English language	Language family
a	a	a	есть (бываетъ)	есть	he is	Суринамский Креольски	Сранан-тонго	Sranan Tongo	Creole
á	a	á	во (въ)	во (в)	in	Романский и древне-Французский	Романский	Romance	Indo-European, Italic, Romance
á	á	á	есть (бываетъ)	есть	he is	Романский и древне-Французский	Романский	Romance	Indo-European, Italic, Romance
a	a	гдѣ	где	where		Романский и древне-Французский	Романский	Romance	Indo-European, Italic, Romance
a	a	во (въ)	во (в)	in		Коптский въ Египтѣ	Коптский в Египте	Coptic in Egypt	Afro-Asiatic, Egyptian
a	a	безъ (кромѣ)	без (кроме)	without		Чукоцкий	Чукотский	Chukchi	Chukotko-Kamchatkan
a	a	онъ	он	he		Суринамский Креольски	Сранан-тонго	Sranan Tongo	Creole
a	a	во (въ)	во (в)	in		Цыганский	Цыганский	Romani	Indo-European, Indo-Aryan
a	a	она	она	she		Суринамский Креольски	Сранан-тонго	Sranan Tongo	Creole
a	a	кѣть, кто	кто	who		Ирландский	Ирландский	Irish	Indo-European, Celtic

First Previous 1 2 3 4 5 ... 6197 Next Last

# The Bantuists' Manifesto

We, the undersigned researchers in African languages, therefore agree:

- to provide our data in electronic form for use by others [...];
- to acknowledge the use we make of the data provided to use by others [...];
- to hold others harmless for any errors or misuse of the data;
- and to make no commercial use of such data [...].

# The Bantuists' Manifesto

A CC-BY-NC license predating the Creative Commons initiative by 5 years - and still a rarity for linguistic online publications like the “Digital Pallas” ...

The screenshot shows a web browser window displaying the homepage of the World Atlas of Language Structures (WALS) Online. The title bar reads "WALS Online - Home". The address bar shows the URL "https://wals.info". The main header features the text "THE WORLD ATLAS OF LANGUAGE STRUCTURES ONLINE" next to a world map where language families are color-coded. Below the header is a green navigation bar with links: Home (selected), Features, Chapters, Languages, References, Authors, Credits, Legal, Download, and Contact. The "Home" link is highlighted with a green background. The main content area has a yellow background and contains the heading "Welcome to WALS Online" and a paragraph describing the project. To the right, there is a blue callout box with text about commenting functionality.

THE WORLD ATLAS  
OF LANGUAGE STRUCTURES  
ONLINE

Home Features Chapters Languages References Authors

Credits Legal Download Contact

## Welcome to WALS Online

The World Atlas of Language Structures (WALS) is a large database of structural (phonological, grammatical, lexical) properties of languages gathered from descriptive materials (such as reference grammars) by a team of 55 authors.

The functionality to comment on WALS Online data via an associated blog has been replaced with mailto links. See also the [contact page](#) for information on how to get in contact with us.

**APiCS**

Main Layout    Main layout    IPA chart    References    Language name    Contributors > SManfredi    Stefano Manfredi    go ▲ set lang    Help  
 Maps    Segments    Examples    Juba Arabic    SPetrollino    Sara Petrollino    go ▼ Print Layout

Feature No. 78 < sort Feature Name Existential verb and transitive possession

By **existential verb** we refer to the element corresponding to English *there is* in existential clauses like *There is food on the table*. This feature asks whether the existential verb is identical to the transitive verb of possession 'have' (cf. value 1 of Feature 77). If your language lacks a transitive verb of possession, choose value 5. Similarly, if your language lacks an existential verb, choose value 5. (The existential verb may of course just be the ordinary copula verb.) See APiCS Glossary ('Identity and differentiation') for a visual representation of the following values:

1/2. The existential verb may be **identical** to the transitive possession verb ('have') (value 1), or it may be a **different** lexical item (value 2).

3/4. If there are two verbs, and one of them means 'there is' and the other 'there is' and 'have', or one of them means 'have' and the other has both meanings, there is **overlap** (value 3). If there are three verbs, and one means only 'have', one means only 'there is', and one means 'there is' and 'have', we have **identity and differentiation** (value 4).

Note that we do not count transitive constructions such as "The table has food on it" as existential constructions -- these are regarded as transitive possession constructions, even though the subject is

Please select only one value words

Values	Value annotation	Value Choice & Confidence > True False	Relative importance	Example request
1 Identity	Haitian <i>gen</i> 'have; exist': Mari <i>gen</i> <i>kouraj</i> 'Mary has courage'	<input checked="" type="radio"/> <input type="radio"/>	% %	164 <i>fì nas bifékir fùu ma</i> go ▲ 165 <i>tiidrá abáo le zol al ma</i> go ▼
2 Differentiation	English <b>have</b> vs. <b>there is</b>	<input type="radio"/> <input checked="" type="radio"/>	% %	Very certain
3 Overlap	Réunion Creole <i>ana</i> 'have; there is' vs. <i>ganye</i> 'have' (in non-present tenses)	<input checked="" type="radio"/> <input type="radio"/>	% %	2
4 Identity and differentiation		<input checked="" type="radio"/> <input type="radio"/>	% %	3
5 The language has no transitive possession verb, or no existential verb	Russian (no transitive possession verb)	<input checked="" type="radio"/> <input type="radio"/>	% %	0

General comments on value assignment

WALS No. [ ] Select if no information available

Feature Source Reference name Pages

Select Remove... Go Own knowledge

Select Remove... Go

More lects add lect  
 Lect my default lect  
 > my default lect view  
 Lect count 1  
 Custom field for database users

The screenshot shows the WALS Online interface. At the top, there's a navigation bar with 'Home', 'Features', and 'Chapters'. Below it, a large section titled 'Chapters' displays a list of entries. The first entry in the list is 'Consonant Inventories', which is highlighted with a light gray background. A modal window is open over this entry, containing the title 'Consonant Inventories' and two tabs: 'Text' (selected) and 'BibTeX'. The 'Text' tab displays the following information:

Ian Maddieson. 2013. Consonant Inventories.  
In: Dryer, Matthew S. & Haspelmath, Martin (eds.)  
WALS Online (v2020.3) [Data set]. Zenodo.  
<https://doi.org/10.5281/zenodo.7385533>  
(Available online at <http://wals.info/chapter/1>, Accessed on 2024-04-02.)

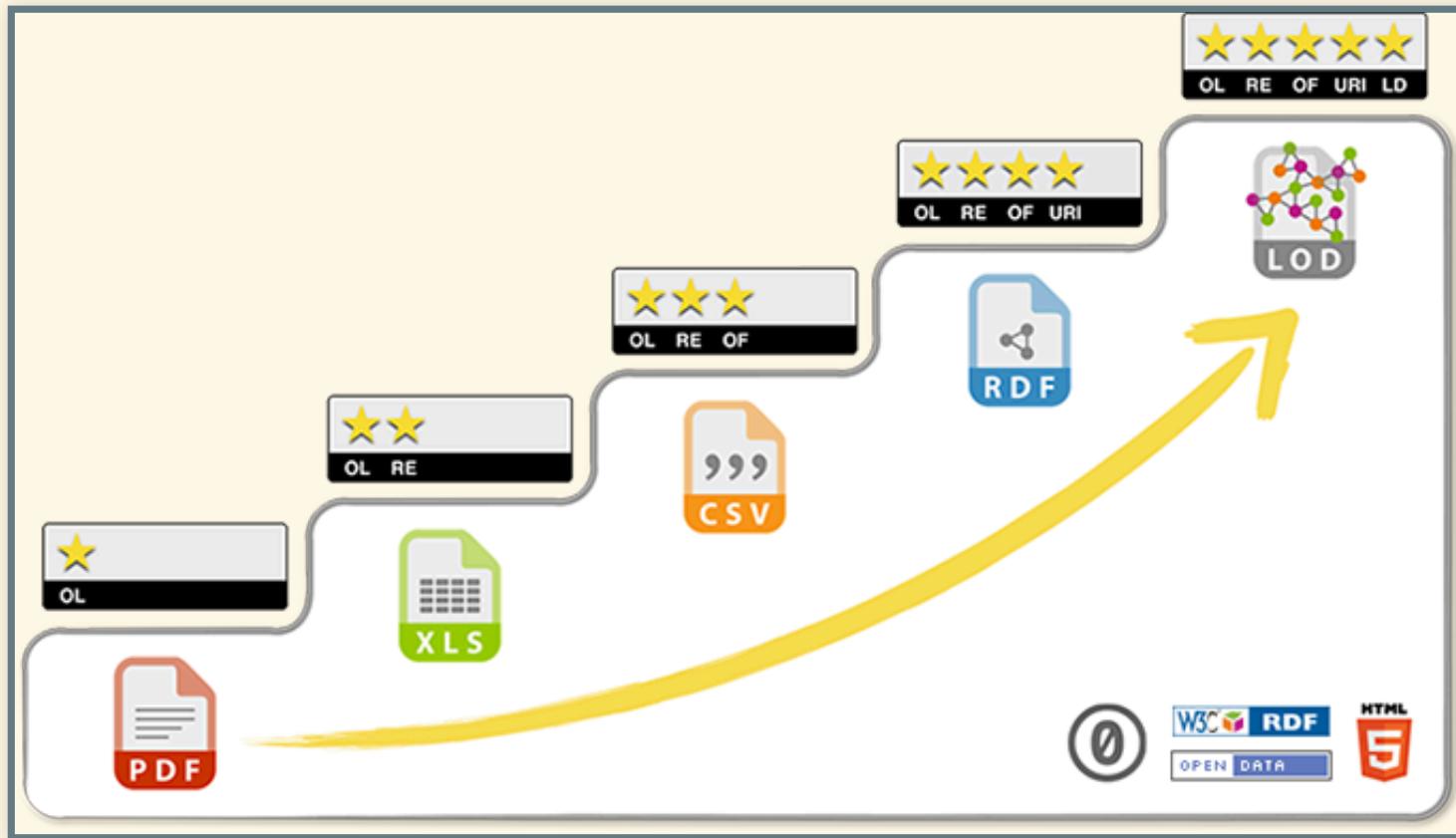
Id	Name	
	<input type="text" value="Search"/>	
	Introduction	
1	Consonant Inventories	Ian Maddieson
2	Vowel Quality Inventories	Ian Maddieson
3	Consonant-Vowel Ratio	Ian Maddieson

# The Heyday of Linguistic Web Apps

## The Semantic Web

*I have a dream for the Web [in which computers] become capable of analyzing all the data on the Web [...]. A “Semantic Web”, which makes this possible, has yet to emerge, but when it does, the day-to-day mechanisms of [...] our daily lives will be handled by machines talking to machines. The “intelligent agents” people have touted for ages will finally materialize*

Tim Berners-Lee, 1998





# Cool URIs don't change

What makes a cool URI?

A cool URI is one which does not change.

What sorts of URI change?

*URIs don't change: people change them.*

# Ruby On Rails In 60 Minutes



RUBY



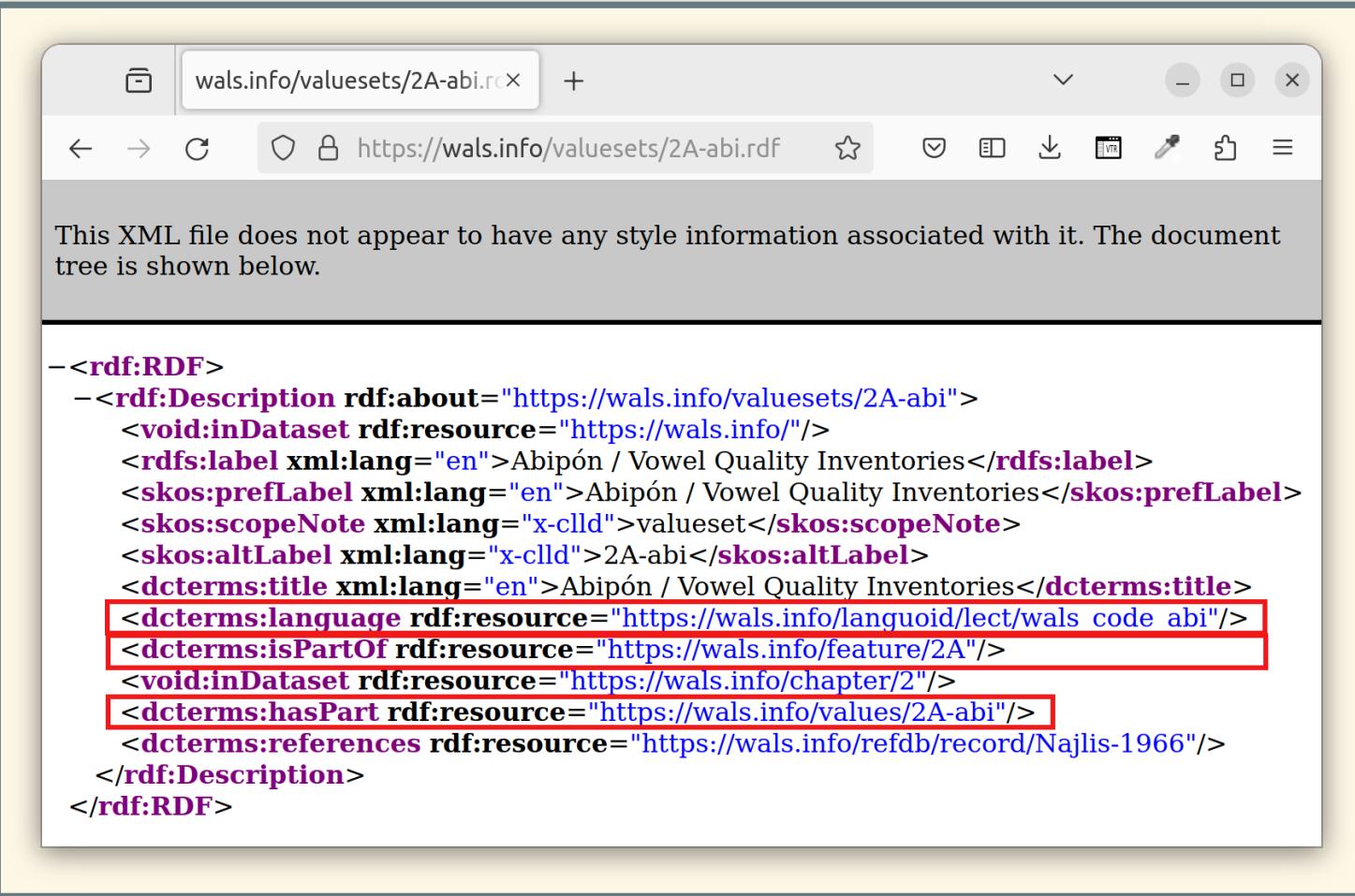
RAILS

60

MIN

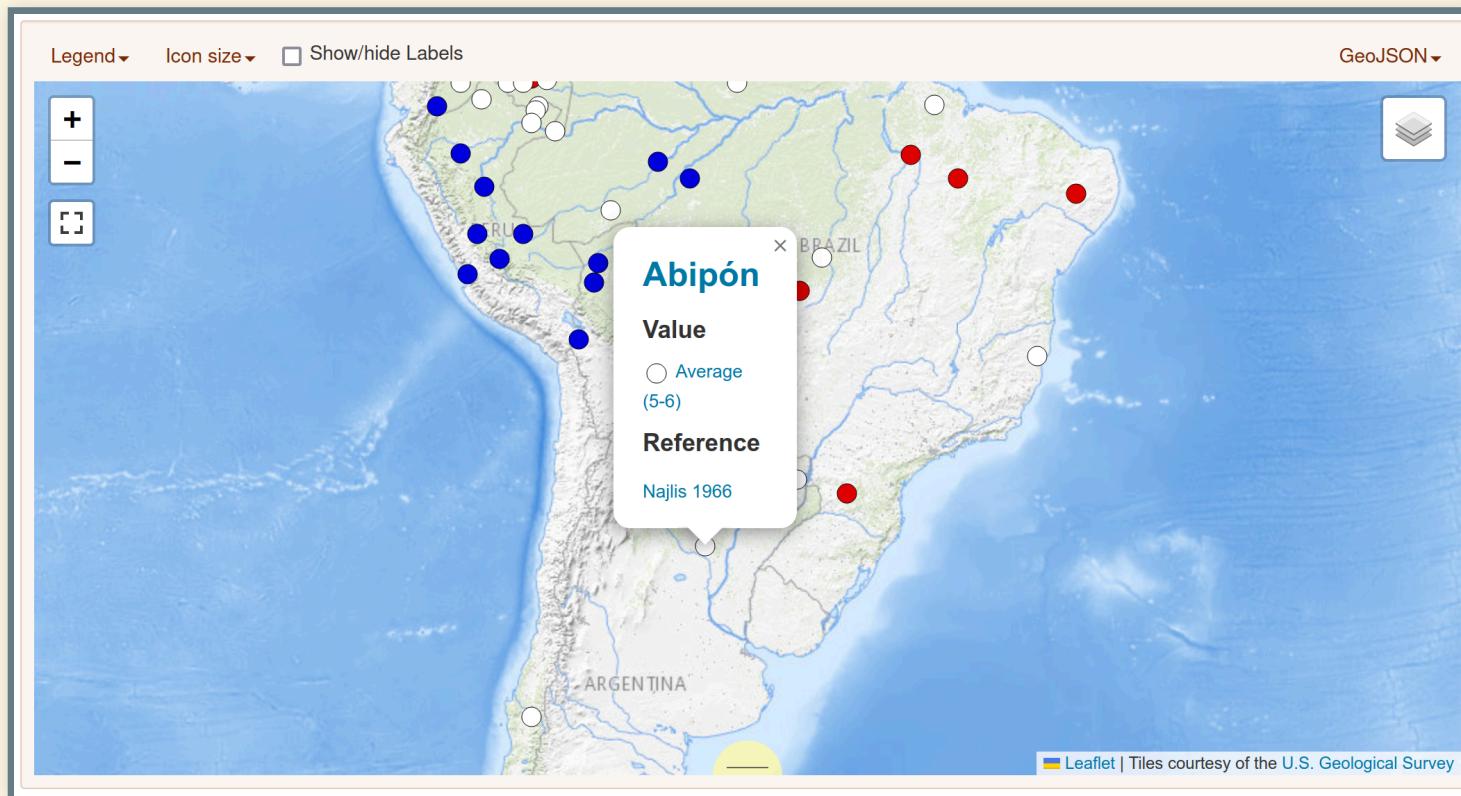
This lead to projects like

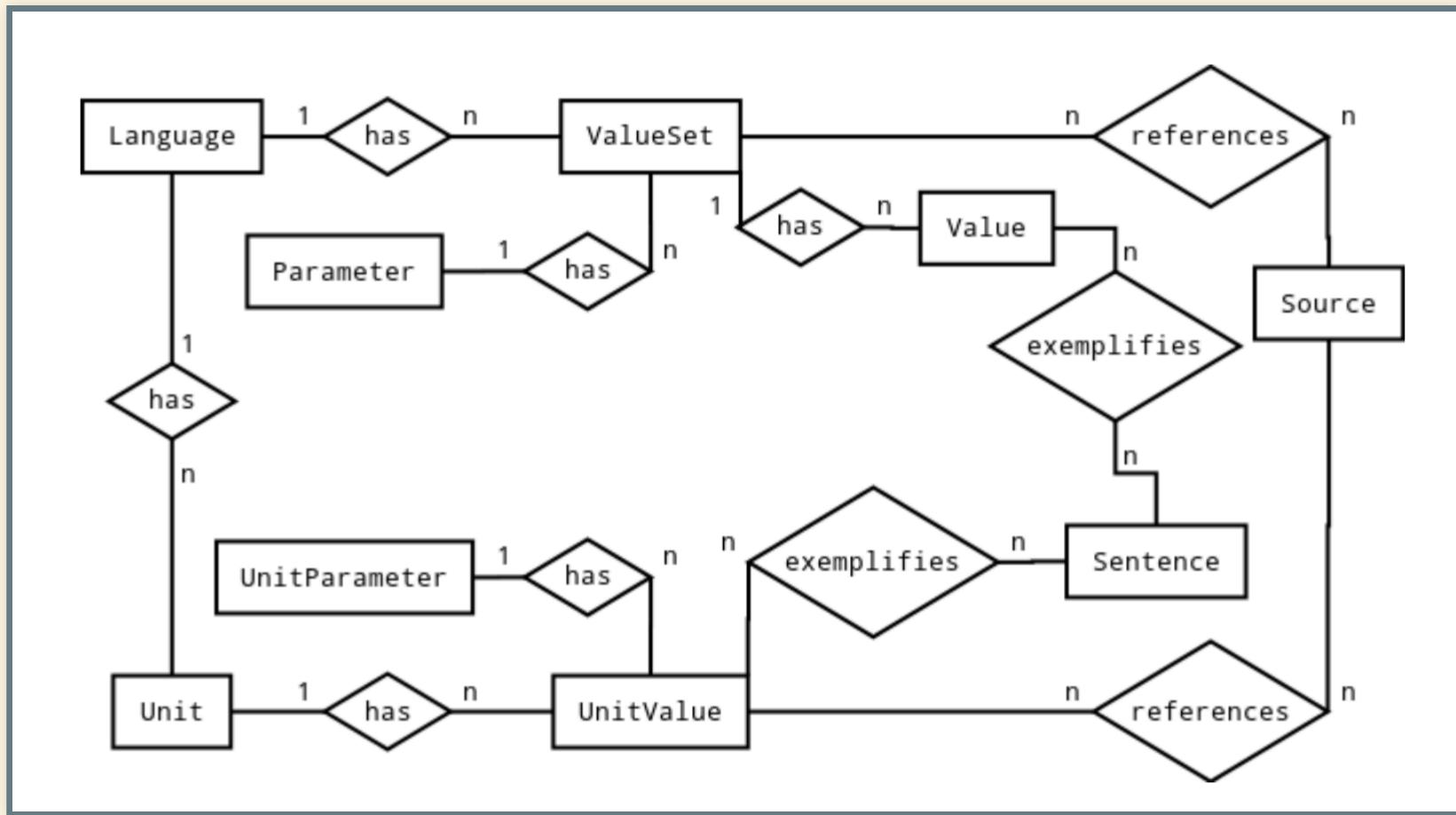
- **GOLD**, the “General Ontology for Linguistic Description” - minting URIs for Linguistic concepts
- **LLOD** - the “Linguistic Linked Open Data” initiative
- **CLLD** - the Cross-Linguistic Linked Data project



This XML file does not appear to have any style information associated with it. The document tree is shown below.

```
--<rdf:RDF>
  --<rdf:Description rdf:about="https://wals.info/valuesets/2A-abi">
    <void:inDataset rdf:resource="https://wals.info/">
    <rdfs:label xml:lang="en">Abipón / Vowel Quality Inventories</rdfs:label>
    <skos:prefLabel xml:lang="en">Abipón / Vowel Quality Inventories</skos:prefLabel>
    <skos:scopeNote xml:lang="x-clld">valueset</skos:scopeNote>
    <skos:altLabel xml:lang="x-clld">2A-abi</skos:altLabel>
    <dcterms:title xml:lang="en">Abipón / Vowel Quality Inventories</dcterms:title>
    <dcterms:language rdf:resource="https://wals.info/languoid/lect/wals_code_abi"/>
    <dcterms:isPartOf rdf:resource="https://wals.info/feature/2A"/>
    <void:inDataset rdf:resource="https://wals.info/chapter/2"/>
    <dcterms:hasPart rdf:resource="https://wals.info/values/2A-abi"/>
    <dcterms:references rdf:resource="https://wals.info/refdb/record/Najlis-1966"/>
  --</rdf:Description>
--</rdf:RDF>
```





- As demonstrated [by the *c11d* framework], a standard software stack is useful.
- But software has a half-life of less than 10 years.
- **Next step is essential: extract a domain specific API which can become standard.**

Forkel & Bank, 2014

# Cracks in the Facade

- Inconsistency
- Versioning
- Maintainability

problems with owl:sameAs



All Images Videos News Shopping More Tools

About 5.920.000 results (0,40 seconds)

In general, the real problem with the use of URIs as identifiers and owl:sameAs is a problem of context and the implicit import of properties. These can all be remedied, and we walk through a case-by-case basis.



W3C

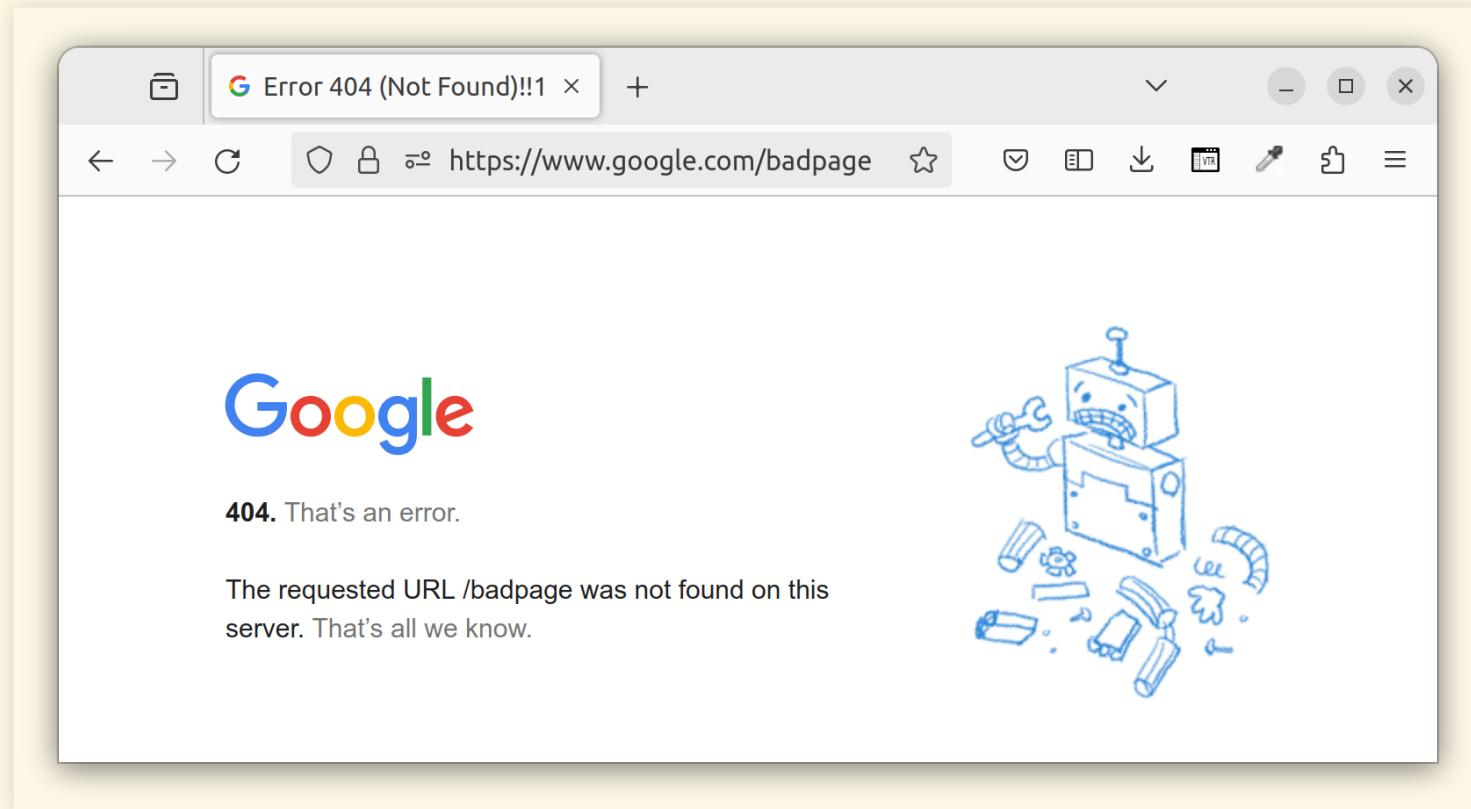
<https://www.w3.org> › 2009/12 › rdf-ws › papers

PDF

:

[When owl:sameAs isn't the Same: An Analysis of Identity Links ...](#)

The screenshot shows a web browser window displaying the WALS Online Changes page. The URL in the address bar is <https://web.archive.org/web/20180404015650/http://wals.info/changes>. The page is part of the Internet Archive Wayback Machine, with 53 captures available. The date shown is April 04, 2018. The main content area features a large world map with language structures. A green navigation bar at the top includes links for Home, Features, Chapters, Languages, References, Authors, Changes (which is active), Credits, Legal, Download, and Contact. To the right, a sidebar titled "Editions" lists the years 2014, 2013, 2011, and 2008.



# Linked Data as persistent API is not maintainable

(Linked Data as data integration paradigm is still useful, though.)

# CLDF - a maintainable API for Linguistic Data

- on-disk format
- line-based text files, suitable for version control (CSV, JSON, BibTeX)
- machine-readable metadata
- Linked Data still in the mix as ontology and links to reference catalogs!

The image shows two side-by-side screenshots. On the left is a GitHub repository interface for the 'wals / cldf /' repository. The 'Files' tab is selected, showing a list of files and a tree view of the 'cldf' folder. A red box highlights the 'cldf' folder and its contents: 'docs', '.gitattributes', 'README.md', 'StructureDataset-metadata.json', 'areas.csv', 'chapters.csv', 'codes.csv', 'contributors.csv', 'countries.csv', 'examples.csv', 'genealogy.csv', 'genealogy.nex', 'language\_names.csv', 'languages.csv', 'media.csv', 'parameters.csv', 'requirements.txt', 'sources.bib', and 'values.csv'. On the right is a screenshot of the WALS dataset's CLDF metadata page. The title is 'StructureDataset The World Atlas of Language Structures Online'. It links to 'CLDF Metadata: StructureDataset-metadata.json' and 'Sources: sources.bib'. Below is a table of CLDF properties and their values, also enclosed in a red box:

property	value
<a href="#">dc:bibliographicCitation</a>	Dryer, Matthew S. & Haspelmath, Martin (eds.) 2013. The World Atlas of Language Structures Online. Leipzig: Max Planck Institute for Evolutionary Anthropology. (Available online at <a href="https://wals.info">https://wals.info</a> )
<a href="#">dc:conformsTo</a>	<a href="#">CLDF StructureDataset</a>
<a href="#">dc:identifier</a>	<a href="https://wals.info">https://wals.info</a>
<a href="#">dc:license</a>	<a href="https://creativecommons.org/licenses/by/4.0/">https://creativecommons.org/licenses/by/4.0/</a>
<a href="#">dcat:accessURL</a>	<a href="https://github.com/cldf-datasets/wals">https://github.com/cldf-datasets/wals</a>
<a href="#">prov:wasDerivedFrom</a>	1. <a href="#">cldf-datasets/wals v2020.2-6-g42c0da7</a> 2. <a href="#">Glottolog v4.6</a>
<a href="#">prov:wasGeneratedBy</a>	1. <b>python:</b> 3.8.10 2. <b>python-packages:</b> <a href="#">requirements.txt</a>
<a href="#">rdf:ID</a>	wals
<a href="#">rdf:type</a>	<a href="http://www.w3.org/ns/dcat#Distribution">http://www.w3.org/ns/dcat#Distribution</a>

## Versions

<a href="#">Version v2020.3</a>	Dec 1, 2022
10.5281/zenodo.7385533	
<a href="#">Version v2020.2</a>	Jul 7, 2022
10.5281/zenodo.6806407	
<a href="#">Version v2020.1</a>	Apr 13, 2021
10.5281/zenodo.4683137	
<a href="#">Version v2020</a>	Mar 27, 2020
10.5281/zenodo.3731125	
<a href="#">Version v2014</a>	Jul 10, 2014
10.5281/zenodo.3607439	

[View all 8 versions](#)

**Cite all versions?** You can cite all versions by using the DOI [10.5281/zenodo.3606197](#). This DOI represents all versions, and will always resolve to the latest one. [Read more](#).

# Another learning process

*should fit in one file*

vs.

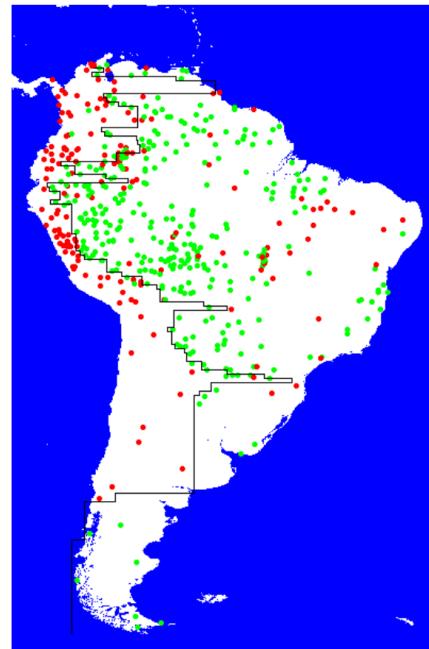
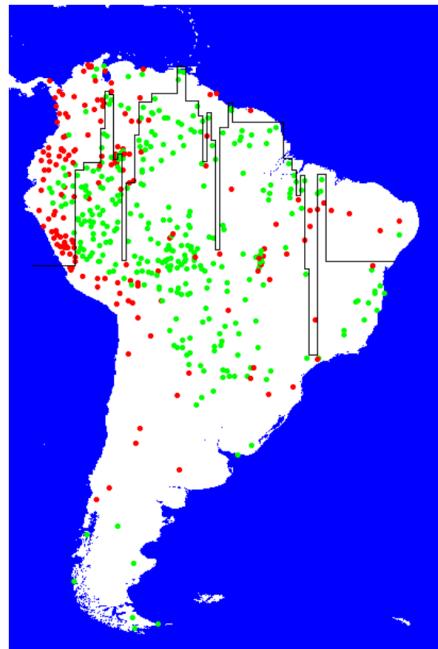
*anticipate the need to use multiple  
tables*

# What does API mean?

- CLDF allows decoupling of dataset curation and analysis tool development!
- Thus, allows replacing the browser as the primary way to interact with data.
- (But `clld` apps are still *one* type of analysis tool for CLDF data).

# Language comparison with CLDF

## Optimal Isogloss Lines: Outcome

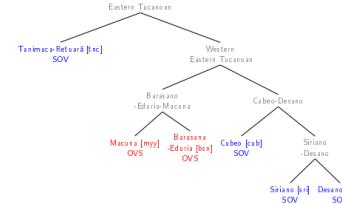


- The optimal line east-west has error 513.125 but
- the optimal line north-south only has 241.875

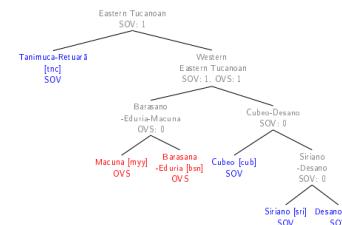
## Parsimony Reconstruct: Procedure & Output

*To each internal node, reconstruct the value that minimizes the total number of changes required*

1. Input (a tree and values at the leaves)



2. For each internal node, starting near the leaves, calculate the minimum number of changes required below it for each possible reconstructed value



3. Reconstruct that which yield the minimum total number of changes in the tree



# Language identification

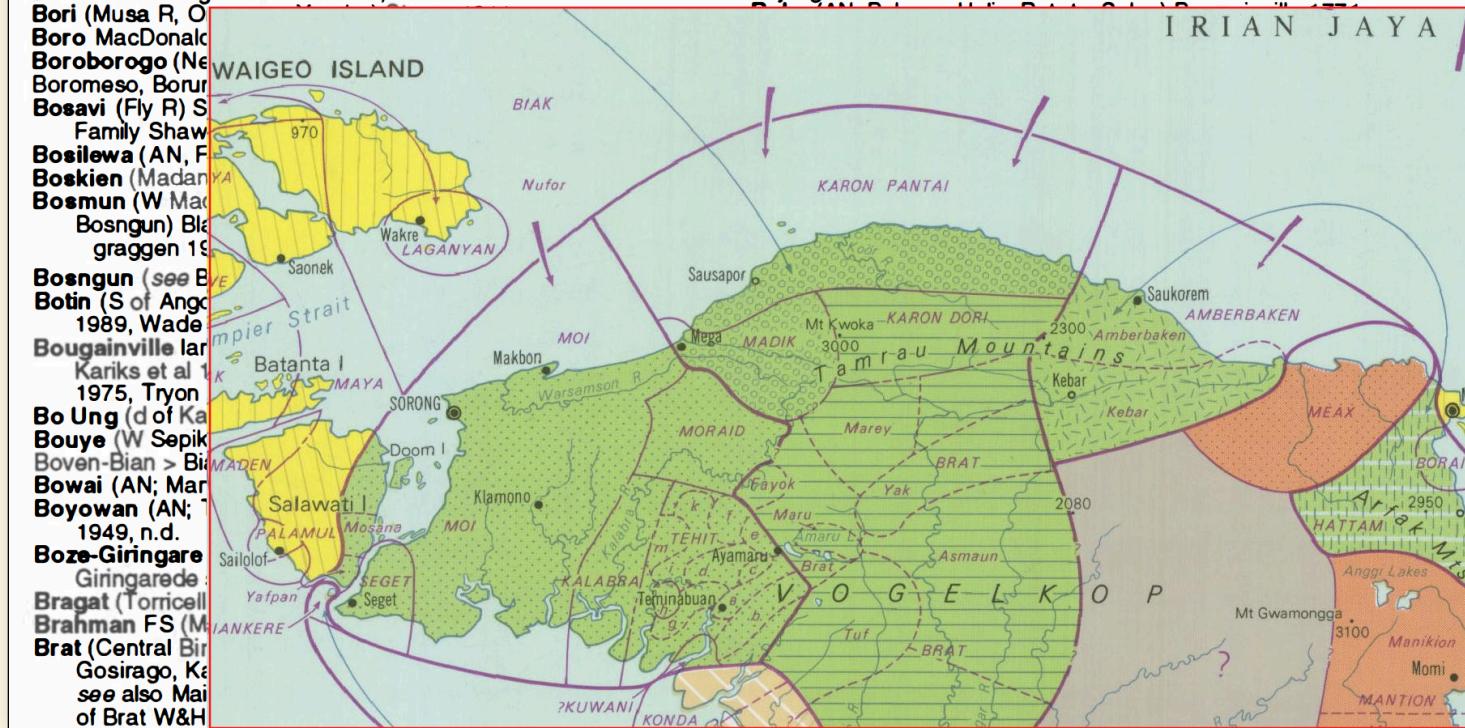
CLDF allows unambiguous identification of languoids  
in a dataset via

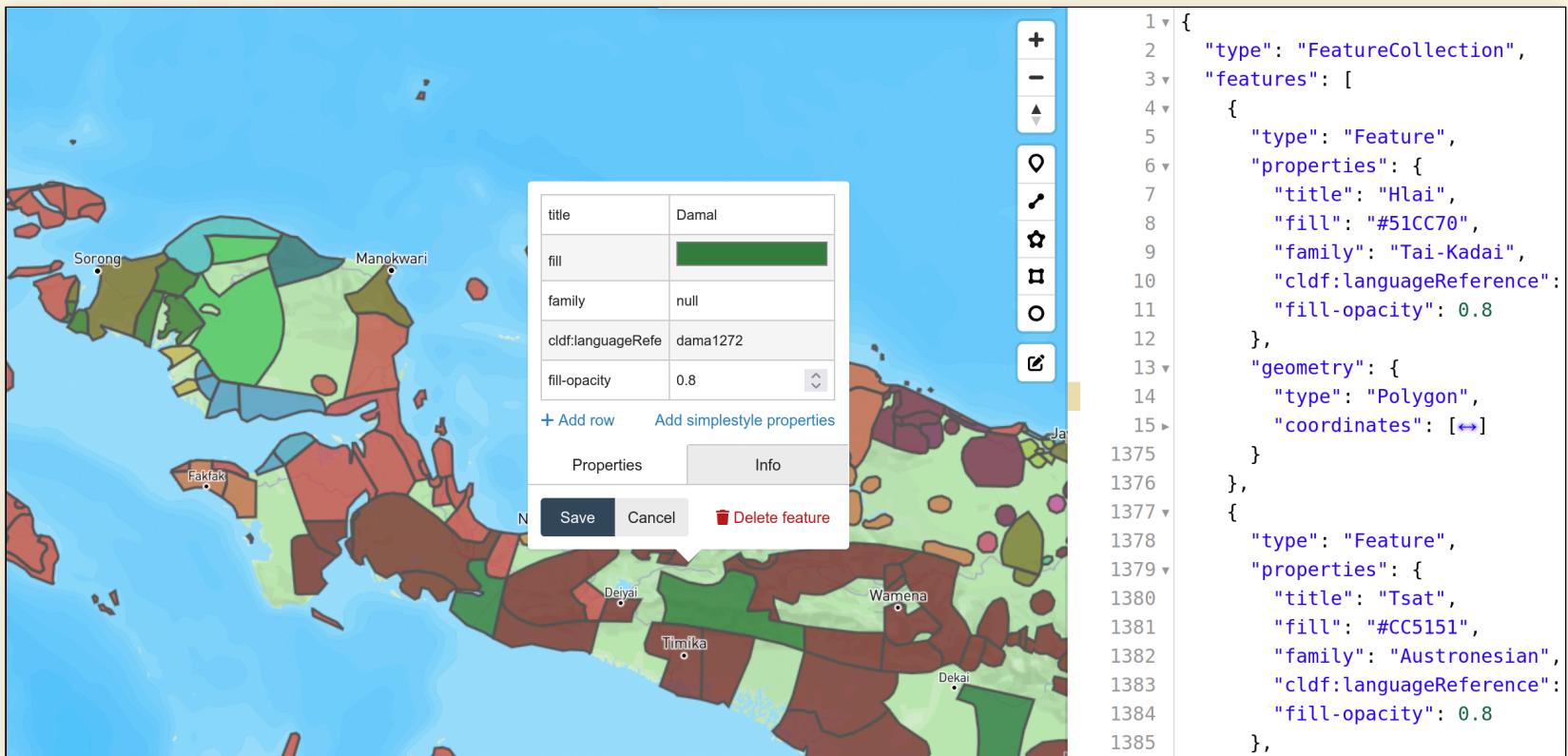
- the LanguageTable component and
- the glottocode property

**Bonkiman** (Madang/Finisterre; also Yupna?) Claassen & McElhanon 1970, McElhanon 1978, W&H 1981, Z'graggen 1975  
**Borai** (E Bird's Head; also Mansim) Barrs 1978, Voorhoeve 1975, W&H 1981; Borai-Hattam W&H 1981  
**Border** d of Komba W&H 1981; Border Stock W&H 1981; Border-Tor-Lake Plain > Northern SS  
**Borebo** d of Magi Thomson 1975, 1975  
**Bori** (Musa R, O)  
**Boro** MacDonald  
**Boroboro** (Ne)  
Boromeso, Borur  
**Bosavi** (Fly R) S Family Shaw  
**Bosilewa** (AN, F)  
**Boskien** (Madan YA)  
**Bosmun** (W Mac)  
Bosngun) Bla graggen 19  
**Bosngun** (see Bo  
**Botin** (S of Ang  
1989, Wade  
**Bougainville** lar  
Kariks et al 1  
1975, Tryon  
**Bo Ung** (d of Ka  
**Bouye** (W Sepik  
Boven-Bian > Bi  
**Bowai** (AN; Mar  
**Boyowan** (AN; T  
1949, n.d.  
**Boze-Giringare**  
Giringaredede  
**Bragat** (Torricell  
**Brahman** FS (M  
**Brat** (Central Bir  
Gosirago, Ka  
see also Mai  
of Brat W&H

Longacre 1972, Lugabai 1971, Müller 1944-45, n.d. x 3,  
Müller & Miltrop 1943-44, n.d., n.d., Nakota 1974, Oliver-Berg  
1979, Potu 1974, Rugabai & Griffin 1971, Schmidt 1909, SIL  
1971, 1975, Simons 1982, Tauria 1974, Thurnwald 1909,  
1912, 1934, 1934, 1937, 1942, n.d., Vaughan 1977, Wheeler  
1911; d's W&H 1981; Buin Family W&H 1981

Buiye > Bouye  
Bujang > Kele





# Concept identification

CLLD Concepticon 3.2.0 - +

https://concepticon.clld.org/contributions

Concepticon Home Concepts Concept sets Concept lists Languages Compilers Sources

## Concept lists

Showing 1 to 2 of 2 entries (filtered from 441 total entries)

Note	Name	Compiler	Alias	Items	Tags	Uniqueness	Year	Gloss languages	Target languages	Sources
	Pallas			Search	Search	--any--	Search	Search	--any--	Search
more	Pallas 1789 285		Pallas, Peter Simon	285	historical	0.04	1789	latin russian	Global	Pallas 1789
more	Pallas 1786 442		Pallas, Peter Simon	442	historical	0.04	1786	french german latin russian	Global	Pallas 1786

Showing 1 to 2 of 2 entries (filtered from 441 total entries)

```
robert@lingn35: ~/projects/concepticon/conceptic... □ ×
```

```
SELECT
    pallas.cldf_name, lsi.cldf_name
FROM
    "concepts.csv" AS pallas,
    "concepts.csv" AS lsi
WHERE
    pallas.cldf_contributionReference = 'Pallas-1789-285' AND
    lsi.cldf_contributionReference = 'Grierson-1928-168' AND
    pallas.cldf_parameterReference = lsi.cldf_parameterReference AND
    lsi.cldf_parameterReference != 0;
```

10,1

All

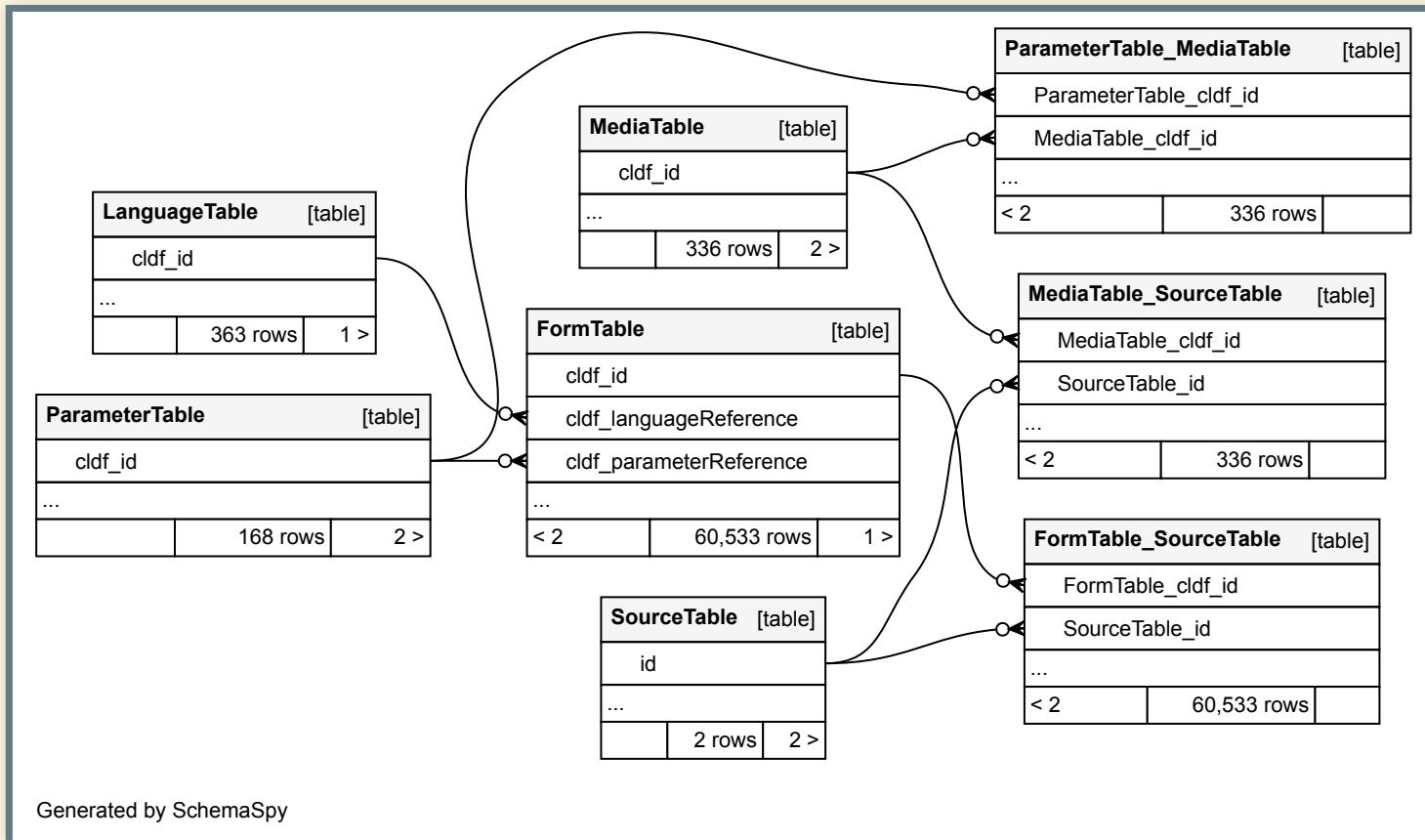
# Exploratory analysis

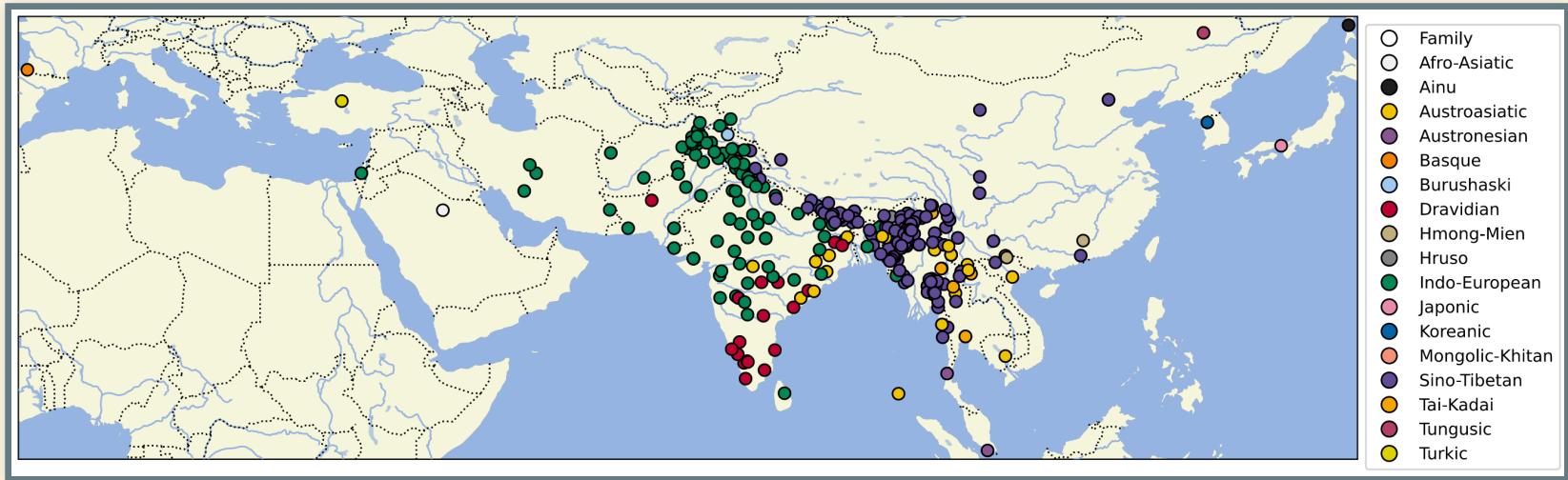
## Table values.csv

property	value
<a href="#">dc:conformsTo</a>	<a href="#">CLDF ValueTable</a>
<a href="#">dc:extent</a>	76475

### Columns

Name/Property	Datatype	Description
<a href="#">ID</a>	string	Primary key
<a href="#">Language_ID</a>	string	References <a href="#">languages.csv::ID</a>
<a href="#">Parameter_ID</a>	string	References <a href="#">parameters.csv::ID</a>
<a href="#">Value</a>	string	
<a href="#">Code_ID</a>	string	References <a href="#">codes.csv::ID</a>
<a href="#">Comment</a>	string	comments in HTML
<a href="#">Source</a>	list of string (separated by ; )	References <a href="#">sources.bib::BibTeX-key</a>
<a href="#">Example_ID</a>	list of string (separated by ; )	References <a href="#">examples.csv::ID</a>





# Analysis methods

## Cognates

CLDF provides all the semantics necessary for LingPy's cognate detection and alignment algorithms:

- unambiguously identified languages and concepts
- segmented word forms

and the data model to store the computed cognacy relations, thus, the input for phylogenetic analyses, e.g. via BEASTling.

# Colexifications

CLDF provides

- the semantics to compute colexifications
- the data model to store the computed networks
- thus, the input for clustering algorithms.

## Generic data access: CLDF SQL

- Each CLDF dataset can be loaded into an SQLite database running a simple command from the `pycldf` package.
- The resulting SQL schema allows uniform data access across datasets.

```
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

# Data modeling

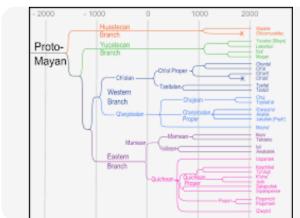
CLDF is extensible and built to evolve by

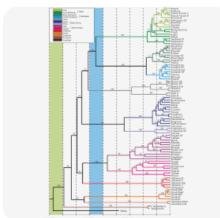
- adding data types (aka tables or components)
- adding properties (aka columns)
- specifying relations between tables/columns and media files

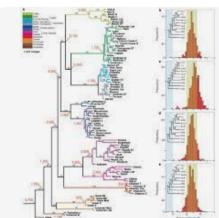
Google language phylogeny

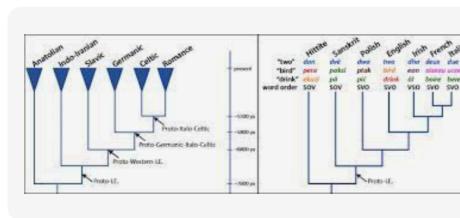
All Images Books News Videos More Tools Saved SafeSearch ▾

english slavic indo-european germanic phylogenetic tree turkic evolution evolution of turkic lang

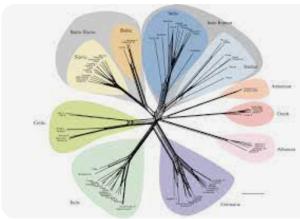
 Tree model - Wikipedia

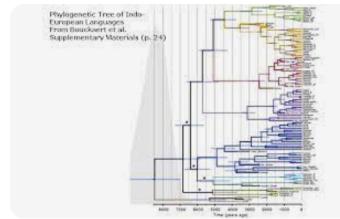
 A dated phylogenetic tree ...

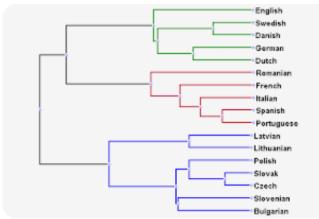
 Language-tree divergence tim...

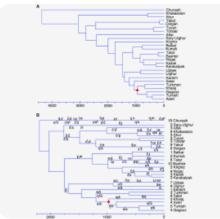
 A phylogeny of the Indo-European ...

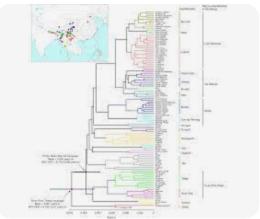
 Comparative linguistics Language evolutionary tre...

 j.g. pausas' blog j.g. pausas' blog x Diversity of languages

 Languages Of The World Malformed Language Tree of Bouckaert ...

 Arya McCarthy Reconstructing Phylogenetic Language ...

 ResearchGate Phylogenetic Trees of the T...

 Nature Phylogenetic evidence for Sino-...

# Aggregation

By providing a uniform way to access data across datasets, CLDF is an essential tool for compiling large, aggregated databases.

# CLDF Wordlist

The screenshot shows a web browser window with the URL <https://lexibank.cldf.org/contributions>. The page title is "Datasets". A note at the top states: "Note that the dataset served here only aggregates a selection of the data from its constituent datasets. Only varieties with at least 100 segmented words for at least 100 different concepts are taken into account. In addition, only one variety per Glottocode - the one with the biggest number of words - from all datasets was included." Below this is a search bar and a table with 8 entries. The table has columns: Id, Name, # languages, # concepts, and # words. A search bar for # languages is set to > 50. The entries are:

Id	Name	# languages	# concepts	# words
blustauronesian	CLDF dataset derived from Blust's Austronesian data coded for the Austronesian Basic Vocabulary Database from 2008	401	210	83,873
transnewguineaorg	CLDF dataset derived from Greenhill's "TransNewGuinea.org" from 2015	253	794	59,787
johanssonsoundsymbolic	CLDF dataset derived from the Johansson et al.'s "The typology of sound symbolism" from 2020	191	284	44,402
bowerpnpy	CLDF dataset derived from Bowern and Atkinson's "Internal Structure of Pama-Nyungan" from 2012	147	338	36,513
polyglottaaficana	CLDF dataset derived from Koelle's "Polyglotta Africana" from 1854	144	262	43,850
northeuralex	CLDF dataset derived from Dellert et al.'s "NorthEraLex (Version 0.9)" from 2020	99	953	105,919
hubercolumbian	CLDF dataset derived from Huber and Reed's "Comparative Vocabulary" from 1992	63	347	23,895
kraftchadic	CLDF dataset derived from Kraft's "Chadic Wordlists" from 1981	58	429	25,076

At the bottom, there is a footer with the lexibank logo, a Creative Commons Attribution 4.0 International License logo, links to Privacy Policy and Disclaimer, a link to the application source on GitHub, and the GitHub logo.

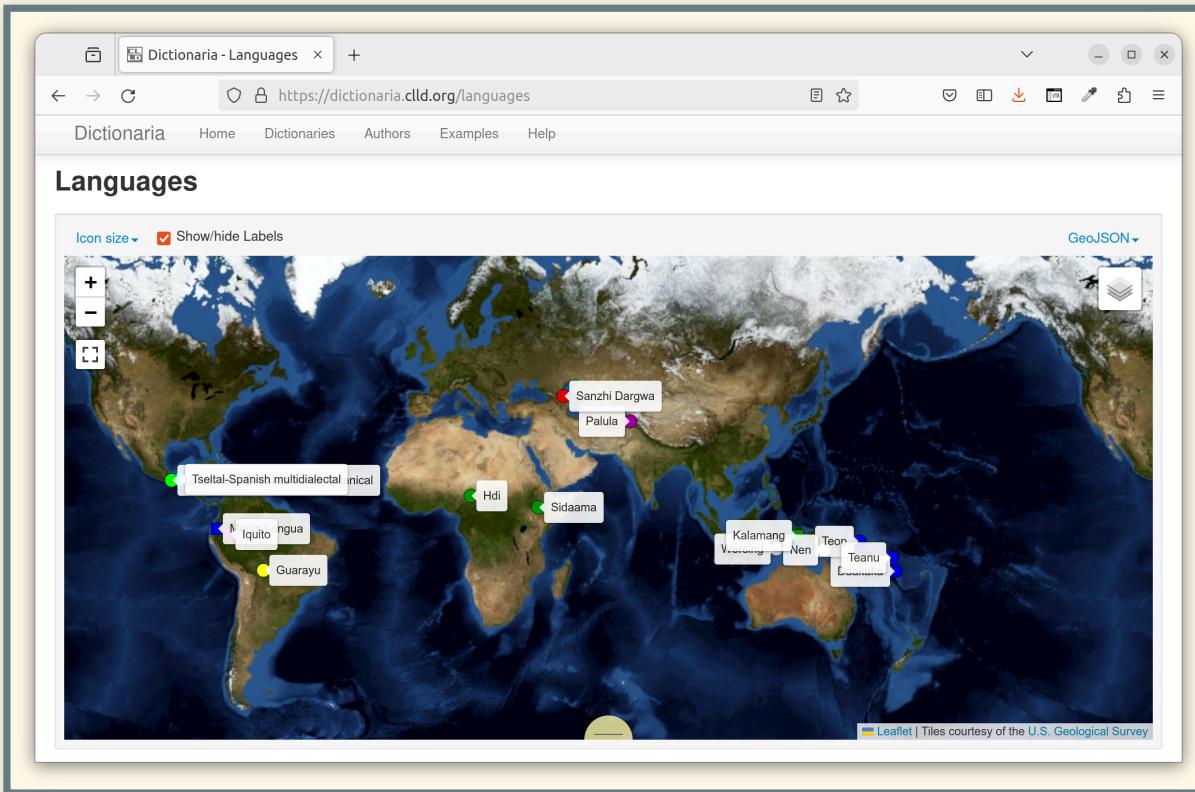
# CLDF StructureDataset

The screenshot shows a web browser window titled "Crossgram - Contribution" at the URL <https://crossgram.cld.org/contributions>. The page has a header with navigation links: Contributions, Languages, L-Parameters, Constructions, Examples, Topics, Sources, and Authors. Below the header is a section titled "Contributions" with the subtext "Showing 1 to 8 of 8 entries". A search bar is present above a table. The table has columns for Name, Authors, Year, and Data source. The data is as follows:

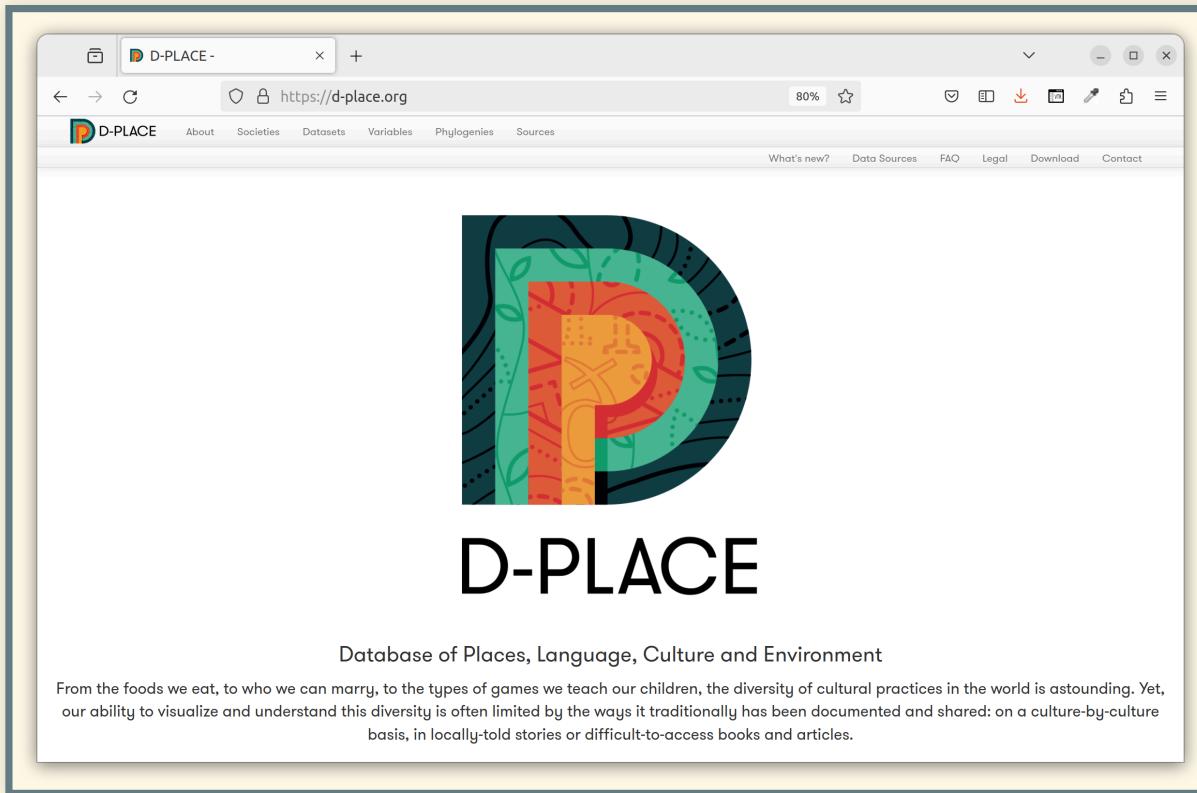
Name	Authors	Year	Data source
Linguistic diversity in space and time	Johanna Nichols	1992	<a href="#">Github: cldf-datasets/nicholsdiversity</a>
Zero marking and the order of core arguments	Kaius Sinnemäki and Noora Ahola	2010	<a href="#">Github: cldf-datasets/sinnemakizeromarking</a>
The 'give' event in Papuan languages	Gerard P. Reesink	2013	<a href="#">Github: cldf-datasets/reesinkgive</a>
Negative existentials: A cross-linguistic study	Ljuba Veselinova	2013	<a href="#">Github: cldf-datasets/veselinovanegex</a>
Order of demonstrative, numeral, adjective, and noun	Matthew S. Dryer	2018	<a href="#">Github: cldf-datasets/dryerorder</a>
Names and nominal classification	Corinna Handschuh	2019	<a href="#">Github: cldf-datasets/handschuhnames</a>
Interrogatives as relativizers in Indo-European	Sandra Auderset	2020	<a href="#">Github: cldf-datasets/audersetinterrog</a>
Estimative constructions cross-linguistically	Guillaume Jacques	2023	<a href="#">Github: cldf-datasets/jacquesestimative</a>

Below the table, there is another search bar and a footer with links to the Max Planck Institute for Evolutionary Anthropology, Leipzig; a Creative Commons Attribution 4.0 International License logo; a Crossgram disclaimer; a GitHub link for the application source; and a GitHub logo.

# CLDF Dictionary



# But also



# More

# CLDF Markdown

## ### Rule 2: Morpheme-by-morpheme correspondence

Segmentable morphemes are separated by hyphens, both in the example and in the gloss. There must be exactly the same number of hyphens in the example and in the gloss. E.g.

[example 2](cldf/examples.csv#cldf:2)

## Rule 2: Morpheme-by-morpheme correspondence

Segmentable morphemes are separated by hyphens, both in the example and in the gloss. There must be exactly the same number of hyphens in the example and in the gloss. E.g.

(2) Lezgian (Haspelmath 1995: 207)

ID	Language_ID	Primary_Text	Analyzed_Word	Gloss	Translated_Text
1					
2	lezg1247	Gila aburun ferma hamışaluğ güğüna amuq'dač.	Gila abur-u-n ferma hamışaluğ güğüna amuq'-dač.	now they-OBL-GEN farm forever behind stay-FUT-NEG	Now their farm will not stay behind forever.

Gila abur-u-n      ferma    hamışaluğ    güğüna    amuq'-da-č.  
now    they-OBL-GEN    farm    forever    behind    stay-FUT-NEG  
'Now their farm will not stay behind forever.'

## Catalog of Linguistic Data

CLDF datasets hold much more “actionable” (meta)data than archive holdings exposed through OLAC.

From

**list all holdings characterized as *text* for language X**

to

**list all segments found in any text for language X**

<https://cldf.clld.org>

