

Guidelines for a standardized data format for use in cross-linguistic studies

October 16 2018



A world map showing data points, for which the researchers plan to gather unified data (e.g., data that is directly comparable) using the guidelines given in the paper. Credit: OpenStreetMap. Forkel et al. 2018. Cross-Linguistic Data Formats, advancing data sharing and re-use in comparative linguistics. *Scientific Data*.

An international team of researchers, members of the Cross-Linguistic Data Formats Initiative (CLDF) led by the Max Planck Institute for the Science of Human History, has proposed new guidelines on crosslinguistic data formats in order to facilitate sharing and data comparisons between the growing number of large linguistic databases worldwide. This format provides a software package, a basic ontology and usage examples.



There is an increasing number of linguistic databases worldwide, raising the possibility of a vast network for potential comparative studies. However, these databases are generally created independently of each other, and often have a unique and narrow focus. This means that the formats used for encoding the data are often different, creating difficulties in comparing data across databases.

The Cross-Linguistic Data Formats Initiative (CLDF) is an effort to resolve these issues. In a paper published in *Scientific Data*, the CLDF sets out proposed guidelines for a standardized format for linguistic databases, and also supplies a software package, a basic ontology and usage examples of best practices. The goal of this effort is to facilitate sharing and re-use of data in comparative linguistics.

The CLDF provides a data model underlying its recommendations that aims to be simple, yet expressive, and is based on the data model previously developed for the Cross-Linguistic Data project. This model has four main entities: (a) languages; (b) parameters; (c) values; and (d) sources. In the model, each value is related to a parameter and a language, and can be based on multiple sources. There are additionally references for sources, and references can also have contexts (which, for example, for printed references would be page numbers).



(a) One Value per Cell	NEITHER:							
		ning	English	German			Dutch	
Many datasets that have been published in the past place	bar	k	bark	Rinde, Borke		ce	bast	
multiple values in the same cell of their data. This is most	NOD	10 D •						
frequently the case with elicitation meanings for which		Meaning English German Dutch						
multiple translations could be found. Since scholars are	hark		bark	Rinde			bact	
rarely explicit about the separators or the techniques by	Dar	A	Dark	Rinde			Dasc	
which they handle these problems, many different ways to	bar	rK	*	BOLKE				
address multiple translations per meaning have been used in	BUT:							
the past, ranging from additional columns up to secondary	ID	Meaning	Language	Form				
shes, pipes) and datasets may even mix the different		bark	English	bark				
techniques. To avoid these problems. CLDF specifies to use	2	bark	German	Rinde			1	
long tables throughout all applications.	3	bark	German	Borke			-	
	4	bark	Dutch	bast			-	
(b) Anticipate the Need of Multiple Tables	NEITHER:							
	Meaning English			German			Dutch	
When a certain complexity of analysis is reached, multiple	bark		A	в, А	, A		с	
tables become inevitable in linguistic datasets.							~	
Unfortunately, the need of multiple tables if often not readily	/ NOR:							
link across tables. Especially formats for cognate coding	Mea	ning	English	German			Dutch	
show great variation in this regard, ranging from multiple	bar	k	bark	Rinde	Borke		bast	
sheets in spreadsheet software that were manually created								
up to customized formats in which additional information is	BUT:							
encoded in form of markup, such as colored cells or text in	TD		-	ADUL A	1	TD	TABLE B	
italic or bold font. All these attempts are very error prone	ID	Meaning	Language	Form		TD	Cognacy	
and lead to data-loss, especially if only certain parts of the	1	bark	English	bark		1	bark-A	
data are snared. To avoid these problems, CLDF specifies to	2	bark	German	Rinde		2	bark-B	
it explicit in the metadate, how tables should be linked	3	bark	German	Borke		3	bark-A	
n explicit in the metadata, now tables should be liftked.	4	bark	Dutch	bast		4	bark-c	

Basic rules of data coding included in the guidelines, taking cognate coding in wordlists as an example. (a) illustrates why long tables should be favored throughout all applications. (b) underlines the importance of anticipating multiple tables along with metadata indicating how they should be linked. Credit: Forkel et al. 2018. Cross-Linguistic Data Formats, advancing data sharing and re-use in comparative linguistics. *Scientific Data*.

The CLDF data model is a package format in which a dataset would be made up of a set of data files containing tables, and a descriptive file that defines the relationships between the tables. Each linguistic data type would have a CLDF module and additional components, which would be the aspects of the data in the module that recur across multiple data



types. The CLDF modules would also contain terms from the CLDF ontology. The ontology is a list of vocabulary that represents objects and properties with well-known semantics in comparative linguistics. This makes it possible for users to reference these terms in a uniform way.

A software package to enable validation and manipulation

The CLDF specifications use common file formats—such as CSV, JSON and BibTeX—that are widely supported, with the goal that these files can easily be read and written on many platforms. Even more importantly, the standardized format will allow researchers without programming skills to access and manipulate the data with preexisting tools, to avoid restricting the package only to researchers with sufficient programming skills to create their own tools. To facilitate this, the CLDF has created a "cookbook" repository for scripts for use with the CLDF specifications.

"We want to bring access to these data and the ability to compare them to as many researchers as possible," says Johann-Mattis List of the Max Planck Institute for the Science of Human History. Robert Forkel, one of the driving forces behind the CLDF initiative, also notes that the CLDF format is not limited to linguistic data alone, but can also incorporate databases of cultural and geographic data, for example. "CLDF may drastically facilitate the testing of questions regarding the interaction between linguistic, cultural, and environmental factors in linguistic and cultural evolution."

More information: Cross-Linguistic Data Formats, advancing data sharing and re-use in comparative linguistics , *Scientific Data*, <u>DOI:</u> <u>10.1038/sdata.2018.205</u>



Provided by Max Planck Society

Citation: Guidelines for a standardized data format for use in cross-linguistic studies (2018, October 16) retrieved 28 April 2024 from <u>https://phys.org/news/2018-10-guidelines-standardized-format-cross-linguistic.html</u>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.