



# Curating Cross-Linguistic Data with git and GitHub

---

Robert Forkel

MEaCoM 1, October 23–24, 2017, Alcanena

Max Planck Institute for the Science of Human History

# CODE IS DATA

Homoiconicity in programming languages means:

*code can be treated as a basic data structure that the programming language knows how to access.*

<http://blogs.mulesoft.org/code-is-data-data-is-code/>

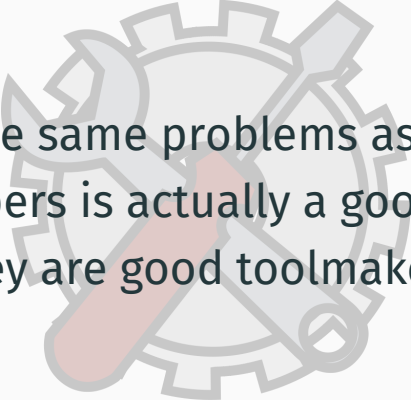
Thus the programming language can manipulate code more reliably (no syntax errors!) than your text editor.

So homoiconicity is a desirable property of a programming language, because it allows for better tooling.

Turns out – in a rather mundane way – data is code, too.

Or less catchy: cross-linguistic data is similar enough to (open source) code to share its tools:

- textual data
- line based
- not really Big Data
- open (often)



Having the same problems as software  
developers is actually a good thing!  
They are good toolmakers!

# DATA CURATION AND VERSION CONTROL

```
commit 1203f919cb811b7fb57fa350f8dc37c0e756a004
Author: kirbykat <kirbykat@users.noreply.github.com>
Date: Thu Oct 5 00:21:47 2017 +0200
```

correct glottolog **id** **for** ancient hebrew

```
commit 3df866f2965a1a9f7e4f67e8dd92833e28baf489
Author: Hans-Jörg Bibiko <bibiko@shh.mpg.de>
Date: Wed Sep 13 09:38:49 2017 +0200
```

fixed Labor **and** Lifecycle spelling

```
commit 8c06135d3fe7cb6f718e48b86e98cdc9e7bb41e9
Merge: 985d3d9 efaf853
Author: xrotwang <xrotwang@googlemail.com>
Date: Wed Sep 13 09:21:31 2017 +0200
```

updated variable definitions

Listing 1: Obviously, data curation and version control software are a perfect match!

## Data curation with GitHub

---

Using line based text files for data and

**git** (a tool for distributed source code management) and

**GitHub** (a platform hosting git repositories).

we get a platform for collaboratively curating cross-linguistic data.

## EXAMPLE: GLOTTOLOG

Open Source software in the age of GitHub is a tremendous success story for worldwide online collaboration.

This is exactly the kind of collaboration we want to enable for data sets like Glottolog, which clearly

- profit from more curators

*given enough eyeballs, all bugs are shallow* (Linus' Law)

- "belong" to the academic community more than to any one institution, thus – given current funding schemes – will have to be transferred to a different owner at some point.



# HOW DO WE TURN GLOTTOLOG INTO OPEN DATA?

We need to model Glottolog data in a way suitable for distributed version control systems.

- line-based text formats, i.e. text that can be meaningfully handled by **diff**
- BibT<sub>E</sub>X for bibliography files
- INI files for languoid metadata.
- A directory tree to model the classification.
- Some tools to simplify manipulation of the language tree.
- An API to access the data in the repository programmatically.

```
@book{94863,  
  address    = {New York},  
  author     = {Sapir, Edward},  
  publisher  = {Harcourt and Brace},  
  title      = {Language},  
  year       = {1949},  
  bibtexkey  = {sapir_language1949},  
  inlg       = {English [eng]},  
  macro_area = {Africa},  
  src        = {wals},  
  srctrickle = {wals#5298}  
}
```

Listing 2: BibT<sub>E</sub>X is used for reference data.

# cldd/glottolog: WHY BibT<sub>E</sub>X?

- Well supported in many bibliography management tools like
  - Zotero
  - jabref
- Our workflow is already adapted to it
- The (missing) details in the data model – e.g. no splitting of authors – align well with our messy data.
- We only use BibT<sub>E</sub>X as container format – no T<sub>E</sub>X in field values, but UTF-8 encoded text.

# cld/glottolog: INI FILES

```
# -*- coding: utf-8 -*-  
[core]  
name = Abinomn  
glottocode = abin1243  
hid = bsa  
level = language  
iso639-3 = bsa  
latitude = -2.92281  
longitude = 138.891  
macroareas =  
    Papunesia  
countries =  
    Indonesia (ID)  
  
[sources]  
glottolog =  
    Mark Donohue and Simon Musgrave 2007 (89329)
```

Listing 3: **INI** files are used for metadata on languoids.

# cldd/glottolog: WHY INI FILES?

- Good support (e.g. syntax highlighting) in many text editors.
- The programming language Python supports reading and writing **INI** files out-of-the-box.
- Format is extensible – new sections and options can be added any time without disrupting the processing pipeline.

```
$ tree --charset ASCII languoids/tree/abkh1242/abkh1243/  
    abkh1244/  
languoids/tree/abkh1242/abkh1243/abkh1244/  
|-- abkh1244.ini  
|-- abzh1238  
|   '-- abzh1238.ini  
|-- bzyb1238  
|   '-- bzyb1238.ini  
'-- samu1242  
    '-- samu1242.ini
```

3 directories, 4 files

Listing 4: A directory tree is used to model the language classification.

## Glottolog and collaboration

---

# THE GITHUB WORKFLOW

**fork** Create your own copy of the data repository. The repository you forked from is also called **upstream**.

**edit** Change the data in your copy.

**commit** Register meaningful groups of changes in your copy.

**pull request** Propose merging your changes into upstream, i.e. **clld/glottolog**.

**merge** Incorporate changes from other forks of the repository.



# USE CASES: TRANSFER OF OWNERSHIP

Forks are essential for the open source software development model for another reason as well:

They allow for seamless transfer of ownership of codebases.

For Glottolog this means

- the data repository can be forked - any fork is as good as the original repos
- the code for the web application has an open license, can be run anywhere, and ingest data from any fork
- the only thing bound to an institution that has to be explicitly transferred (with consent of the owner) is the domain name **glottolog.org**

# USE CASES: FUNCTIONALITY BUILT ON THE REPOSITORY

Functionality can be built on top of the repository – rather than on top of the web application

- reduces traffic at **glottolog.org**
- works off-line
- works for forks, too, ...
- ...thus, local changes can be incorporated in workflows right away
- put an API on your data rather than on your web apps

# GLOTTOLOG REFERENCE SEARCH

```
$ glottolog refsearch "author:Holton_year:2003_Tobelo"
ID                               Author          Year  Title
-----
mpieva:Holton2003Tobelo         Holton, Gary    2003  Tobelo
wals:2737                       Holton, Gary    2003  Tobelo
hh:g:Holton:Tobelo             Gary Holton     2003  Tobelo
langsci:Holton:03              Holton, Gary    2003  Tobelo
(4 matches)
```

Listing 5: Using the pure-python Whoosh search engine, we can provide full-text search on more than 300,000 Glottolog references from the command line (or any program that can "shell out")

# GLOTTOLOG LANGUOID SEARCH

```
$ glottolog langsearch Deutsch
4 matches
German Sign Language [germ1281] language
./sign1238/deaf1237/dgsi1234/germ1281/md.ini
Deutsche Gebärdensprache

Kaniet-Dempwolff [kani1283] language
./aust1307/nucl1752/.../west2532/anch1239/kani1283/md.ini
hh:hw:Dempwolff:Deutsch-Neuguinea** recorded

Old Saxon [olds1250] language
./indo1319/germ1287/.../alts1234/olds1250/md.ini
Altnieder-deutsch

German [stan1295] language
./indo1319/germ1287/.../high1287/stan1295/md.ini
Deutsch
```

Listing 6: Similar functionality is provided to search language information.

# USE CASES: ADD "YOUR" LANGUAGE

Working on varieties which are not in Glottolog?

- mint Glottocodes (using functionality built on top of the repository)
- add languoids to your fork of the repository
- use "your" Glottocodes in your data ...
- ...while waiting for "upstream" to incorporate your changes.


## OTHER USEFUL LINE BASED TEXT FORMATS


- bagit: cataloging/packaging hierarchies of files
- csv: tabular data for version control
- csv packages with w3c: multi-table packages with foreign keys
- cldf: cross-linguistic datatypes built on the w3c spec for tabular data


```
myfirstbag/  
|-- data  
|   |-- 27613-h  
|       |-- images  
|           |-- q172.png  
|           |-- q172.txt  
|-- manifest-md5.txt  
|   49afbd86a1ca9f34b677a3f09655eae9 data/27613-h/images/q172  
|   .png  
|   408ad21d50cef31da4df6d9ed81b01a7 data/27613-h/images/q172  
|   .txt  
|-- bagit.txt  
    BagIt-Version: 0.97  
    Tag-File-Character-Encoding: UTF-8
```


Listing 7: <https://en.wikipedia.org/wiki/BagIt>

**correct glottolog id for ancient hebrew** [Browse files](#)

 master

 kirbykat committed 6 days ago 1 parent af0c317 commit 1203f919cb811b7fb57fa350f8dc37c0e756a004

 Showing 1 changed file with 1 addition and 1 deletion. Unified Split

2  datasets/SCCS/societies.csv View ▼


	@@ -53,7 +53,7 @@ SCCS57,xd576,Kurd,cent1972,Kurd (SCCS57),,1951,Lang_assignment_change_notes,Kurd
53	53 SCCS55,xd577,Abkhaz,abkh1244,Abkhaz (SCCS55),,1880,Lang_assignment_change_notes,Abkhaz (RI03),http://ehrafworldcultures.
54	54 SCCS47,xd581,Turks,nuc11301,Turks (SCCS47),,1950,Lang_assignment_change_notes,Turks (MB01),http://ehrafworldcultures.
55	55 SCCS46,xd588,Rwala,east2690,Rwala Bedouin (SCCS46),Rwala Bedouin,1913,Lang_assignment_change_notes,Rwala (MD04),http://
56	-SCCS44,xd589,Hebrews,hebr1245,Hebrews (SCCS44),,-621,Lang_assignment_change_notes,,,31.18,34.92,31.18,34.92,Original
56	+SCCS44,xd589,Hebrews,anci1244,Hebrews (SCCS44),,-621,Lang_assignment_change_notes,,,31.18,34.92,31.18,34.92,Original
57	57 SCCS45,xd590,Babylonians,akka1240,Babylonians (SCCS45),,-1750,Lang_assignment_change_notes,,,32.58,44.75,32.58,44.75,
58	58 SCCS59,xd600,Punjabi,panj1256,Punjabi (West) (SCCS59),Punjabi (West) ,1950,Lang_assignment_change_notes,,,32.5,74.0,3
59	59 SCCS58,xd605,Basser1,bass1257,Basser1 (SCCS58),,1958,Lang_assignment_change_notes,Basser1 (MA10),http://ehrafworldcul

Figure 1: CSV plays well with version control and GitHub.



## Bells and whistles

---

Let's go further borrowing best practices in software development.

## *Continuous integration*

*In addition to automated [...] tests, organisations using CI typically use a build server to implement continuous processes of applying quality control in general — small pieces of effort, applied frequently.*

*[http://en.wikipedia.org/wiki/Continuous\\_integration](http://en.wikipedia.org/wiki/Continuous_integration)*

# CI FOR GITHUB

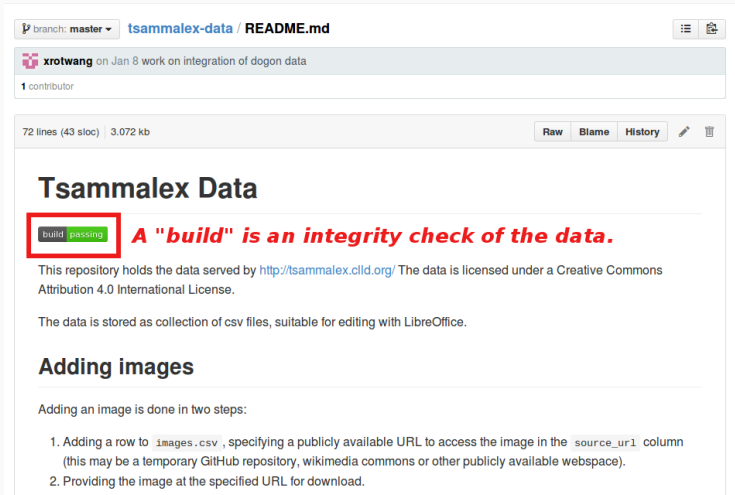


Figure 2: GitHub repositories can be registered with CI service provider Travis-CI.

# CI BUILD HISTORY

cld/tsammalex-data build failing

<https://travis-ci.org/cld/tsammalex-data>

Current Branches **Build History** Pull Requests Settings

	master fixed problems	# 213 passed	1 min 6 sec
	xrotwang committed	525ed3e	about 11 hours ago
	master names ids correction	# 212 failed	1 min 9 sec
	LenaSell committed	afbc572	about 13 hours ago
	master gwj categories and habitats added	# 211 failed	1 min 9 sec
	LenaSell committed	6dbe981	about 13 hours ago
	master Comparative data from Bantu languages (c	# 210 failed	1 min 31 sec
	ChristfriedNaumann committed	6d2927f	about 15 hours ago
	master updated taxa info from external sources	# 209 passed	1 min 53 sec
	xrotwang committed	6cfc023	about 19 hours ago
	master Bibliography extended	# 208 passed	1 min 15 sec
	ChristfriedNaumann committed	574ba97	a day ago

Figure 3: The build history relates builds and repository changes.

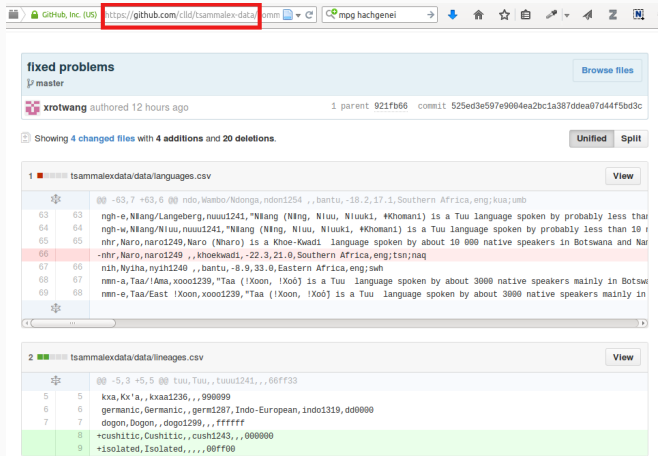
# CI BUILD LOG

*The build log identifies errors and allows line-specific linking.*

```
4965     self.test(*self.arg)
4966     File "/home/travis/build/tsamalex-data/tsamalexdata/tests/test_csv.py"
4967     raise ValueError('integrity checks failed!')
4968 nose.proxy.ValueError: integrity checks failed!
4969 ----- >> begin captured stdout << -----
4970 ERROR:languages:66: non-unique id: nhr
4971 ERROR:names:5120: invalid reference viljoenkamupingene1983[41]
4972 ERROR:names:5121: invalid reference viljoenkamupingene1983[41]
4973 ERROR:languages:6: invalid lineages id referenced: cushitic
4974 ERROR:languages:12: invalid lineages id referenced: cushitic
4975 ERROR:languages:23: invalid lineages id referenced: isolated
4976 ERROR:languages:26: invalid lineages id referenced: cushitic
4977 ERROR:languages:81: invalid lineages id referenced: isolated
4978
4979 ----- >> end captured stdout << -----
4980
4981 -----
4982 Ran 1 test in 2.955s
4983
4984 FAILED (errors=1)
4985
4986 The command "nosetests" exited with 1.
4987
4988 Done. Your build exited with 1.
```

Figure 4: Build log for error reporting.

# CI: ADDRESSING BUILD ERRORS



The screenshot shows a GitHub commit page for repository `tsammaxlex-data`. The commit message is `fixed problems` by user `xrotwang`, authored 12 hours ago. The commit hash is `525ed3e597e9004ea2bc1a387ddea07d44f5bd3c`. The commit log shows 4 changed files with 4 additions and 20 deletions. The files are:

- `tsammaxlexdata/data/languages.csv` (View)
- `tsammaxlexdata/data/lineages.csv` (View)

The `languages.csv` file shows a diff with the following changes:

Line	Old	New
63	ngh-e,Nlång/Langeberg,nuuu1241,"Nlång (Nlång, Nluu, Nluuki, #Khomani) is a Tuu language spoken by probably less than 10 native speakers in Botswana and Namibia"	ngh-e,Nlång/Langeberg,nuuu1241,"Nlång (Nlång, Nluu, Nluuki, #Khomani) is a Tuu language spoken by probably less than 10 native speakers in Botswana and Namibia"
64	ngh-w,Nlång/Nluu,nuuu1241,"Nlång (Nlång, Nluu, Nluuki, #Khomani) is a Tuu language spoken by probably less than 10 native speakers in Botswana and Namibia"	ngh-w,Nlång/Nluu,nuuu1241,"Nlång (Nlång, Nluu, Nluuki, #Khomani) is a Tuu language spoken by probably less than 10 native speakers in Botswana and Namibia"
65	nhr,Naro,naro1249,"Naro (Naro, Nluu, Nluuki, #Khomani) is a Tuu language spoken by about 10 native speakers in Botswana and Namibia"	nhr,Naro,naro1249,"Naro (Naro, Nluu, Nluuki, #Khomani) is a Tuu language spoken by about 10 native speakers in Botswana and Namibia"
66	nhr,Naro,naro1249,,khoekwadi,-22.3,21.0,Southern Africa,eng;tsn;naq	nhr,Naro,naro1249,,khoekwadi,-22.3,21.0,Southern Africa,eng;tsn;naq
67	nih,Nyiha,nyih1240,,bantu,-8.9,33.0,Eastern Africa,eng;sw	nih,Nyiha,nyih1240,,bantu,-8.9,33.0,Eastern Africa,eng;sw
68	nmn-a,Taa/I Ama,xoo1239,"Taa (I Xoon, I Xoo) is a Tuu language spoken by about 3000 native speakers mainly in Botswana and Namibia"	nmn-a,Taa/I Ama,xoo1239,"Taa (I Xoon, I Xoo) is a Tuu language spoken by about 3000 native speakers mainly in Botswana and Namibia"
69	nmn-e,Taa/East I Xoon,xoo1239,"Taa (I Xoon, I Xoo) is a Tuu language spoken by about 3000 native speakers mainly in Botswana and Namibia"	nmn-e,Taa/East I Xoon,xoo1239,"Taa (I Xoon, I Xoo) is a Tuu language spoken by about 3000 native speakers mainly in Botswana and Namibia"

The `lineages.csv` file shows a diff with the following changes:

Line	Old	New
5	kka,Kx'a,,kkaa1236,,999999	kka,Kx'a,,kkaa1236,,999999
6	germanic,Germanic,,germ1287,Indo-European,indo1319,dd0000	germanic,Germanic,,germ1287,Indo-European,indo1319,dd0000
7	dogon,Dogon,,dого1299,,ffffff	dogon,Dogon,,dого1299,,ffffff
8	+cushitic,Cushitic,,cush1243,,000000	+cushitic,Cushitic,,cush1243,,000000
9	+isolated,Isolated,,,,00ff00	+isolated,Isolated,,,,00ff00

Figure 5: The URL to the build log could be used in a commit log to link changes back to the error report.

# BUT WHAT IF GITHUB ...?

Does this introduce too much dependence on GitHub.com?

There are some mitigating factors:

- git is a *distributed scm*, thus each clone contains all the data!
- There are alternative git hosting platforms like BitBucket.
- and then there's ZENODO

ZENODO solves the longterm preservation and citability issue for GitHub repositories by

- archiving releases ("issues") of GitHub repositories
- assigning a DOI to each release



# GLOTTOLOG 2.4 AT ZENODO

The screenshot shows the Zenodo interface for the dataset 'glottolog-data: Glottolog database 2.4'. The header includes the Zenodo logo and the tagline 'Research. Shared.'. Navigation links for Search, Communities, Browse, Upload, and Get started are present, along with Sign In and Sign Up buttons. The dataset is dated 20 March 2015 and is marked as 'Software' and 'Open access'. The authors listed are Harald Hammarström, Robert Forkel, Martin Haspelmath, and Sebastian Bank. A red box highlights the citation text: 'Hammarström, Harald & Forkel, Robert & Haspelmath, Martin & Bank, Sebastian. 2015. Glottolog 2.4. Leipzig: Max Planck Institute for Evolutionary Anthropology. (Available online at <http://glottolog.org>)'. Below this is a table of files with columns for Name, Date, and Size. A red box highlights the file 'glottolog-data-v2.4.zip' with a date of '20 Mar 2015' and a size of '377.4 MB', and a 'Download' button. To the right, a 'GitHub' badge indicates the dataset is available on GitHub. Further down, the 'Publication date' is confirmed as '20 March 2015', and a red box highlights the 'DOI' as '10.5281/zenodo.16245'. The 'Keyword(s)' listed is 'linguistics'. The 'Related publications and datasets' section includes a link to the GitHub repository: 'https://github.com/cld/glottolog-data/tree/v2.4'.

zenodo Research. Shared.

Search Communities Browse Upload Get started Sign In Sign Up

20 March 2015 Software Open access

glottolog-data: Glottolog database 2.4

Harald Hammarström; Robert Forkel; Martin Haspelmath; Sebastian Bank

(show affiliations)

Hammarström, Harald & Forkel, Robert & Haspelmath, Martin & Bank, Sebastian. 2015. Glottolog 2.4. Leipzig: Max Planck Institute for Evolutionary Anthropology. (Available online at <http://glottolog.org>)

Name	Date	Size
glottolog-data-v2.4.zip	20 Mar 2015	377.4 MB

Download

Available in

GitHub

Publication date: 20 March 2015

DOI: 10.5281/zenodo.16245

Keyword(s): linguistics

Related publications and datasets:

Supplement to: <https://github.com/cld/glottolog-data/tree/v2.4>

Figure 6: <http://dx.doi.org/10.5281/zenodo.16245>

If your data is code, treat it as such.

And yes, GitHub is the missing editorial backend of your system.

