

# Glottocodes: Identifiers Linking Families, Languages and Dialects to Comprehensive Reference Information

Robert Forkel<sup>a,\*</sup> and Harald Hammarström<sup>b</sup>

<sup>a</sup> *Department of Linguistic and Cultural Evolution, Max Planck Institute for Evolutionary Anthropology, Germany*

*E-mail: robert\_forkel@eva.mpg.de*

<sup>b</sup> *Department of Linguistics and Philology, Uppsala University, Sweden*

*E-mail: harald.hammarstrom@lingfil.uu.se*

**Abstract.** Glottocodes constitute the backbone identification system for the language, dialect and family inventory Glottolog (<https://glottolog.org>). In this paper, we summarize the motivation and history behind the system of glottocodes and describe the principles and practices of data curation, technical infrastructure and update/version-tracking systematics. Since our understanding of the target domain — the dialects, languages and language families of the entire world — is continually evolving, changes and updates are relatively common. The resulting data is assessed in terms of the FAIR (Findable, Accessible, Interoperable, Reusable) Guiding Principles for scientific data management and stewardship. As such the glottocode-system responds to an important challenge in the realm of Linguistic Linked Data with numerous NLP applications.

Keywords: Linguistics, Linked Data, Language Inventory, Linguistic Standards

## 1. Introduction

Glottocodes constitute the backbone identification system for the language, dialect and family inventory Glottolog (<https://glottolog.org>, currently in edition 4.4, [1]). A glottocode consists of four alphanumeric characters (i.e., lowercase letters or decimal digits) and four decimal digits, for example `abcd1234` or `b10b1234`. Glottocodes are complementary to three-letter ISO 639-3 language identification codes (see <https://iso639-3.sil.org/>) which, however, concern languages only.

In the current release, there are 25,900 glottocodes (8,533 language-level, 4,571 family-level and 12,796 dialect-level).

---

\*Corresponding author. E-mail:  
[robert\\_forkel@eva.mpg.de](mailto:robert_forkel@eva.mpg.de).

## 2. Motivation and History

Glottocodes were introduced in 2010 by Glottolog collaborator Sebastian Nordhoff, in response to the following requirements:

- An ID specifically designed for machine readability, not confusable with an informal or human-directed identifier
- An ID type oblivious to level of linguistic abstraction (idiolect, sociolect, dialect, language, subfamily, family, etc.)
- An ID system for languages that improves on the ISO 639-3 language identifiers in terms of quality, transparency and anchoring

The 8-character long alphanumeric string was designed not to resemble an abbreviation or an easily remembered mnemonic. This was done specifically in order to counter any temptation to

1 capitalize, modify, inflect or translate it, which  
 2 users might if the ID-string had had a more  
 3 human-palatable appearance (such as a three-  
 4 letter mnemonic, a standardized name or the like).

5 Glottolog had adopted a doculect-based ap-  
 6 proach<sup>1</sup> for organizing concrete attestations of lan-  
 7 guages as recorded by bibliographical references.  
 8 This means that language data (ultimately ema-  
 9 nating from idiolects of specific speakers) recorded  
 10 in different publications are grouped into succes-  
 11 sively larger conglomerates such as subdialects, di-  
 12 alects, languages, subfamilies and families [7]. Ear-  
 13 lier approaches had always sought to tie the iden-  
 14 tity of a set of data to a specific level such as  
 15 language or dialect, such that if the set of data  
 16 did not change but the level changed (e.g., a lan-  
 17 guage reconsidered a dialect or vice versa) the ID  
 18 had to change as well (e.g., [8]). Given controver-  
 19 sies over language and dialect status, along with  
 20 our incomplete understanding of the language situ-  
 21 ation for many minority languages around the  
 22 world, such level changes are actually not un-  
 23 common. With the level-independent ID-system  
 24 of glottocodes, merely a level-attribute, but not  
 25 the ID itself would have to change in such cases.  
 26 The level-neutral term for the denotation of a gлот-  
 27 tocode is *languoid* [2].

28 Finally, a decade or longer ago, the quality and  
 29 transparency of the ISO 639-3 standard for lan-  
 30 guages was problematic and an alternative was  
 31 clearly needed [9]. ISO 639-3<sup>2</sup> aims at complete  
 32 coverage of all natural languages and thus super-  
 33 sedes the earlier ISO 639-2 and ISO 639-1 stan-  
 34 dards. SIL International is the registration author-  
 35 ity for ISO 639-3 and also publish the Ethnologue  
 36 language catalogue [10]. ISO 639-3 carries little  
 37 metadata and/or justification for its entries as part  
 38 of the standard but the information in the corre-  
 39 sponding entries in Ethnologue is in practice of-  
 40 ten taken to substitute as such. Fortunately, the  
 41 quality and transparency has improved in the last  
 42 decade so that the discrepancy between language-  
 43 level glottocodes and ISO 639-3 codes is dimin-  
 44 ishing (but not completely eradicated, see Section  
 45 3.3).  
 46  
 47

48 <sup>1</sup>The term 'doculect' emanates from [2] but the general  
 49 approach has been used by many earlier authors, notably  
 50 [3–6].

51 <sup>2</sup>See <https://iso639-3.sil.org/>.

1 Glottolog was initiated by Harald Hammarström,  
 2 Sebastian Nordhoff and Martin Haspelmath and is  
 3 now run by a group of editors<sup>3</sup>. Editors are chosen  
 4 per release based on unanimous approval of the  
 5 previous group of editors.  
 6  
 7

### 8 3. Glottolog data curation

9  
 10 As explained above, glottocodes are the identi-  
 11 fiers for *languoids* — the main objects of the data  
 12 curated in Glottolog. Glottolog also curates bibli-  
 13 ographical references ('langdoc') in much the same  
 14 way, though this is not the focus of the present  
 15 paper.

16 In the following, we briefly describe the infras-  
 17 tructure and framework pioneered for the Gлот-  
 18 tolog data for data curation and publication ([11]  
 19 and [12]).

- 20 – All data is stored in UTF-8 encoded text files  
 21 (with consistency ensured by the `pyglottolog`  
 22 software package ([13])).
- 23 – Thus, collaboration and curation workflows  
 24 can make use of `git`, a distributed version  
 25 control system to track history of changes and  
 26 provenance.
- 27 – The master copy of the Glottolog `git` reposi-  
 28 tory is hosted with GitHub at <https://github.com/glottolog/glottolog> for curation (but due  
 29 to the nature of git repositories — where all  
 30 essential metadata is part of every copy of the  
 31 repository — this does not put the data at  
 32 the mercy of GitHub).
- 33 – Released versions of the repository are pub-  
 34 lished and archived with Zenodo (<https://zenodo.org>).

35  
 36  
 37  
 38 This setup does not only provide a stable man-  
 39 agement system for all information shared between  
 40 the Glottolog editors, but also a collaborative envi-  
 41 ronment which allows involving the wider commu-  
 42 nity. Somewhat similar to — but less formal than  
 43 — ISO 639-3 change requests, Glottolog users can  
 44 make use of GitHub issues to indicate errors or re-  
 45 quest inclusion of new languoids<sup>4</sup>, or even submit  
 46 pull requests with data corrections<sup>5</sup>.  
 47

48 <sup>3</sup>[https://github.com/glottolog/glottolog/blob/master/](https://github.com/glottolog/glottolog/blob/master/CONTRIBUTORS.md)  
 49 `CONTRIBUTORS.md`

50 <sup>4</sup>E.g. <https://github.com/glottolog/glottolog/issues/646>

51 <sup>5</sup>E.g. <https://github.com/glottolog/glottolog/pull/648>

1 Glottolog aims to share this data in an open  
2 and FAIR ([14]) way. Stepping through the FAIR  
3 Guiding Principles for scientific data management  
4 and stewardship<sup>6</sup> will shed light on the details.

### 6 3.1. Glottolog data is findable

8 Glottolog is a well established language catalog  
9 as evidenced by more than 400 citations of edi-  
10 tions of Glottolog such as "Glottolog 4.0" in the  
11 scholarly literature (according to Google Scholar).  
12 The Glottolog data repository lists 16 contribu-  
13 tors in addition to the Glottolog editors (and not  
14 including users opening issues, see <https://github.com/glottolog/glottolog/graphs/contributors>) —  
15 pointing to a healthy, collaborative user commu-  
16 nity.

17 *Glottolog data are registered and indexed in a*  
18 *searchable resource.*

21 Glottolog data is also well indexed in relevant  
22 catalogues: The first point of contact for many  
23 users is the Glottolog web application at <https://glottolog.org> — not at least because it is well  
24 indexed by Google and other search engines. But  
25 Glottolog data is also harvested by OLAC (Open  
26 Language Archives Community) and is listed in  
27 the OLAC catalogue as the archive with the high-  
28 est number of distinct languages (see <http://www.language-archives.org/metrics/glottolog.org>). Fi-  
29 nally, Glottolog data releases can be found on Zen-  
30 odo, and consequently wherever Zenodo metadata  
31 is indexed.

32 *Glottolog data are assigned a globally unique*  
33 *and persistent identifier.*

34 All languoids in Glottolog are unambiguously  
35 identified via glottocodes. These glottocodes are  
36 transparently associated with URLs in the glottolog.  
37 org domain, turning them into globally unique  
38 identifiers. Each release of Glottolog is identified  
39 by the DOI assigned by Zenodo.

40 *Metadata clearly and explicitly include the*  
41 *identifier of the data they describe.*

42 CLDF — one of the dissemination formats of  
43 Glottolog — is designed to allow for explicit link-  
44 ing of metadata to identifiers. The underlying

1 mechanism to do this is described in [15], and the  
2 semantics are provided through the CLDF Ontol-  
3 ogy (see <https://cldf.cldf.org/v1.0/terms.rdf>).

### 5 3.2. Glottolog data is accessible

6 *Glottolog data are retrievable by their identifier*  
7 *using a standardised communications protocol.*

9 Zenodo (and the metadata associated by Zen-  
10 odo with the DOI assigned to data releases) guar-  
11 antees that data is retrievable using the standard  
12 protocol associated with DOIs.

13 For each languoid, the CLDF/CSVW data as-  
14 sociates an HTTP ([16]) URL, which is resolvable  
15 via the Glottolog web application.

### 17 3.3. Glottolog data is interoperable

18 *Glottolog data use a formal, accessible, shared,*  
19 *and broadly applicable language for knowledge*  
20 *representation.*

21 Glottolog aims at integration with the Seman-  
22 tic Web at large and the Linguistic Linked Data  
23 initiative in particular.

24 At the most fundamental level this means re-  
25 source URLs — aka URLs for Glottocodes. These  
26 resource URLs are not only usable as universally  
27 unique identifiers, but are also resolvable through  
28 the Glottolog web application. HTTP status codes  
29 ([16]) returned by the web application signal the  
30 status of Glottocodes as follows:

31 **200 OK** for active codes

32 **410 GONE** or **301 MOVED PERMANENTLY** for re-  
33 tired Glottocodes

34 **404 NOT FOUND** for invalid codes

35 The web application also provides several serial-  
36 izations of RDF ([17]) representations of languoid  
37 data. These serializations can be retrieved using  
38 standard content negotiation mechanisms such as  
39 using ACCEPT HTTP headers.

40 While these efforts provide convenient integra-  
41 tion with the "living" Semantic Web, Glottolog  
42 also aims at interoperability for its archived, long-  
43 term available datasets. To this end, Glottolog  
44 data is serialized as a CLDF Structure Dataset  
45 ([18]). The CLDF standard ([19], [20]) does not  
46 only provide interoperability with other CLDF  
47 datasets, but — due to being built on the W3C's  
48 "CSV on the Web" recommendation ([15] and [21])  
49 — also allows automatic conversion to RDF ([22]).

50 <sup>6</sup><https://www.go-fair.org/fair-principles/>

1 *Glottolog data use vocabularies that follow*  
2 *FAIR principles.*

3 CLDF bundles data with structured, machine  
4 readable, semantic web-ready metadata. Since  
5 CLDF metadata is encoded in JSON-LD ([23]),  
6 the data can be marked up using standard ontolo-  
7 gies such as Dublin Core, DCAT ([https://www.](https://www.w3.org/ns/dcat#)  
8 [w3.org/ns/dcat#](https://www.w3.org/ns/dcat#)) and PROV ([https://www.w3.](https://www.w3.org/ns/prov#)  
9 [org/ns/prov#](https://www.w3.org/ns/prov#)).

10 *Glottolog data include qualified references to*  
11 *other (meta)data.*

12 Thanks to improvements in the curation of ISO  
13 639-3 language identifiers during the last decade,  
14 ISO 639-3 codes and language-level glottocodes  
15 are one-to-one interchangeable for the vast major-  
16 ity of cases, and the differences are few enough  
17 that a specific comment explaining the differences  
18 are given in each of the remaining cases on Glot-  
19 tolog. In fact, Glottolog aims at covering all valid  
20 ISO 639-3 codes to provide a full mapping, but  
21 typically there is a time lag of a couple of months  
22 between additions to ISO 639-3 and a Glottolog  
23 release addressing these changes. There remains a  
24 principled difference in anchoring where the deno-  
25 tation of a glottocode in Glottolog is defined by the  
26 data and information in the references tied to it.  
27 The references are associated in Glottolog in such  
28 a way that the referenced data and information is  
29 enough to distinguish the languoid from all other  
30 languoids. Strictly speaking, the ISO 639-3 stan-  
31 dard provide no definition or justification of the  
32 recorded entries. In Ethnologue [10] — the refer-  
33 ence for most of the ISO 639-3 codes — each entry  
34 has metadata such as geographical information,  
35 name(s), speaker numbers and classification which  
36 presumably defines the language, but no actual or  
37 referenced data from the denoted language. Unfor-  
38 tunately, it is not so that metadata information is  
39 in all cases enough to identify its denotation. Lan-  
40 guage names are notoriously ambiguous and the  
41 case of language-shifting ethnic groups is particu-  
42 larly tricky, as most metadata (speaker numbers,  
43 geography, name) is not sufficient to disambiguate  
44 between the original and substituted language.

45 Glottolog, and in particular individual lan-  
46 guoids, are also well-linked from Wikipedia. In  
47 particular, practically all language- and family-  
48 level languoids are referenced in Wikipedia. These  
49 Wikipedia links translate to Wikidata links (e.g.  
50  
51

1 <https://www.wikidata.org/wiki/Q31746>) which in  
2 turn provide links to other language identification  
3 schemes such as ISO 639-2 (which are arguably less  
4 important than ISO 639-3 in the contexts where  
5 Glottolog is most used).  
6

### 7 3.4. *Glottolog data is reusable*

8 *Glottolog data are released with a clear and ac-*  
9 *cessible data usage license.*

10 Glottolog data is release under a CreativeCom-  
11 mons CC-BY-4.0 license.  
12  
13

14 *Glottolog data are associated with detailed*  
15 *provenance.*

16 Like most large-scale databases, parts of Glot-  
17 tolog data are aggregated from various sources.  
18 Glottolog tries to be transparent about this,  
19 e.g. by  
20

- 21 – providing references for all classification propos-  
22 als<sup>7</sup>
- 23 – providing references for all endangerment as-  
24 sessments<sup>8</sup>
- 25 – describing the provenance of the bibliogra-  
26 phy<sup>9</sup>

27 *Glottolog data meet domain-relevant commu-*  
28 *nity standards.*

29 We already described the relation between glot-  
30 tocodes and ISO 639-3 language codes. Arguably,  
31 the transparent mapping between the two, which  
32 Glottolog provides, is the most important con-  
33 tribution towards meeting domain-relevant stan-  
34 dards.  
35

36 But as explained above, Glottolog also  
37

- 38 – caters to the LLD community, by meeting Se-  
39 mantic Web standards, e.g. re-using ontolo-  
40 gies like GOLD<sup>10</sup> to identify languoid levels  
41 and Lexvo.org<sup>11</sup> to identify ISO 639-3 codes  
42 in Glottolog's RDF formats,  
43

44 <sup>7</sup>E.g. [https://github.com/glottolog/glottolog/blob/v4.](https://github.com/glottolog/glottolog/blob/v4.4/languoids/tree/indo1319/anat1257/luvi1234/cari1274/md.ini#L59-L63)  
45 [4.4/languoids/tree/indo1319/anat1257/luvi1234/cari1274/](https://github.com/glottolog/glottolog/blob/v4.4/languoids/tree/indo1319/anat1257/luvi1234/cari1274/md.ini#L59-L63)  
46 [md.ini#L59-L63](https://github.com/glottolog/glottolog/blob/v4.4/languoids/tree/indo1319/anat1257/luvi1234/cari1274/md.ini#L59-L63)

47 <sup>8</sup>See [https://github.com/glottolog/glottolog/blob/v4.4/](https://github.com/glottolog/glottolog/blob/v4.4/config/aes_sources.ini)  
48 [config/aes\\_sources.ini](https://github.com/glottolog/glottolog/blob/v4.4/config/aes_sources.ini)

49 <sup>9</sup>See [https://github.com/glottolog/glottolog/blob/v4.4/](https://github.com/glottolog/glottolog/blob/v4.4/references/BIBFILES.ini)  
50 [references/BIBFILES.ini](https://github.com/glottolog/glottolog/blob/v4.4/references/BIBFILES.ini)

51 <sup>10</sup><http://purl.org/linguistics/gold/>

<sup>11</sup><http://lexvo.org/>

- 1 – serves the OLAC community, by implement-
- 2 ing the OAI-PMH data provider specification
- 3 ([24]), thereby allowing harvesting through
- 4 OLAC,
- 5 – helps researchers in descriptive and compar-
- 6 ative linguistics to inform their analyses us-
- 7 ing Glottolog metadata, by making this data
- 8 accessible as CLDF dataset,
- 9 – provides the NLP community with the means
- 10 necessary to follow the "Bender Rule" [25]
- 11 of always identifying the language(s) (or lan-
- 12 guage varieties) involved in NLP research

#### 14 4. Policies governing glottocode

#### 15 assignment

16  
17  
18 Glottolog aims to be complete with respect to  
19 all assertable L1 languages<sup>12</sup> in the real world, so  
20 all languages in the world (as far as this is un-  
21 derstood at the time of a certain release) have a  
22 language-level glottocode. Glottolog makes a clas-  
23 sification decision for all language-level languoids  
24 so the family-level inventory is complete in the  
25 sense of exhausting the languages of a given re-  
26 lease.

27 Glottolog also classifies dialects insofar as it  
28 attaches them to exactly one language-level lan-  
29 guoid. But the inventory of dialects (varieties of a  
30 language), non-L1 languages (artificial languages,  
31 speech registers, pidgins) and non-assertable lan-  
32 guages (putative languages for which there is in-  
33 sufficient data to decide if they are different from  
34 all other languages) and putative families (hy-  
35 potheses about family relationships that have ap-  
36 peared in the literature) is growing but still far  
37 from complete. The world may contain more of any  
38 or all of these entities without a necessary reflec-  
39 tion in a glottocode. Genuine completeness with  
40 respect to these categories is deemed practically  
41 (if not theoretically) impossible.

42 For these reasons, Glottolog accounts for any  
43 changes to the language-level inventory between  
44 two releases, i.e. language-level glottocodes of the  
45 previous release will always be valid glottocodes in  
46 the next. So if something was deemed a real-world  
47 language, a user can follow any changes to that  
48 assertion. If a language-level languoid was com-

49  
50 <sup>12</sup>See <https://glottolog.org/glottolog/>  
51 glottologinformation for an explanation of these criteria.

1 pletely erroneous, it is moved to the Bookkeep-  
2 ing category. If it is promoted/demoted to a fam-  
3 ily/dialect, it retains its glottocode but changes its  
4 level accordingly, but from that point on it ceases  
5 to be “protected” by its language-level status, so  
6 may be retired in the next release<sup>13</sup>. In contrast,  
7 the family-level and dialect-level glottocodes are  
8 not “protected” and may be removed from the in-  
9 ventory between releases. Since they do not nec-  
10 essarily reflect a real-world entity like an L1 lan-  
11 guage, it cannot systematically be explained what  
12 “happened” to them, e.g., if they never really ex-  
13 isted. However, some tracking possibilities are al-  
14 ways guaranteed because both families and di-  
15 alects are linked to language-level languoids. If a  
16 family-level glottocode disappears, it is possible  
17 to check which language-level languoids it covered  
18 and to check which family/ies they are now asso-  
19 ciated with and if a dialect-level glottocode disap-  
20 pears, it is possible to check which language-level  
21 languoid it pertained to and to check which di-  
22 alects are now associated with it. Furthermore, the  
23 git-versioning and the structure of the index allows  
24 a quick location of the specific pull-request associ-  
25 ated with the removal/appearance of a glottocode  
26 (see Section 5).

27 Glottocodes are not recycled — for new enti-  
28 ties, completely new glottocodes are assigned (re-  
29 tired codes are not re-used/re-purposed). Hence,  
30 all glottocodes that have ever appeared are either  
31 active or retired.

#### 32 5. Glottolog versioning

33  
34  
35 We already pointed out that Glottolog data is  
36 versioned and released periodically (aiming at a  
37 bi-annual release frequency). Each such Glottolog  
38 release is self-contained, i.e. does not reference any  
39 base data, but instead includes it. Thus, when link-  
40 ing other resources to Glottolog, one should always  
41 specify the particular target version.

42 Glottolog follows a semantic-versioning scheme<sup>14</sup>  
43 for data:

44  
45  
46  
47 <sup>13</sup>This process works similar to the way depreca-  
48 tion (<https://en.wikipedia.org/wiki/Deprecation>) is used  
49 in software development to provide limited backwards com-  
50 patibility.

51 <sup>14</sup><https://semver.org/>

- 1 – Resources using Glottolog should always target the highest patch version of a particular minor version. This should not break any processing code, but may correct errata.
- 2
- 3
- 4
- 5 – Upgrading resources to a new minor version may change data/links, but should not break processing code.
- 6
- 7
- 8 – Upgrading to a new major version may break processing code, i.e. the data structure may change.
- 9
- 10

11 The Glottolog version history can be explored in two ways: The Glottolog web application resolves resource URLs of obsolete languoids as follows:

```
12 $ curl -I https://glottolog.org/resource/languoid/id/awun1244
13 HTTP/1.1 301 Moved Permanently
14 Date: Fri, 08 Jan 2021 12:06:32 GMT
15 Content-Type: text/html; charset=UTF-8
16 Location: https://glottolog.org/files/glottolog-4.0/awun1244.html
```

17 Where the HTML page at <https://glottolog.org/files/glottolog-4.0/awun1244.html> provides context about the languoid and versions when it was still active. In the case of obsolete dialect-level languoids — such as Awuna — this typically allows determining the parent language-level languoid, which will always be part of the current release.

21 Alternatively, since Glottolog data is curated as git repository, we can use the git software to inspect the history. In the case of Awuna (awun1244), e.g., we learn that it was removed during a "cleanup" of the classification of the Gbe sub-family:

```
22 $ git log --all --full-history -- "**/awun1244/md.ini"
23 commit 923cd9eb21f27bc9ae0797a315dd48fd009e53c4
24 Author: d97hah <harald@bombo.se>
25 Date: Thu Jul 25 14:50:11 2019 +0200
26 Gbe (#387)
27 * Gbe resolved clf + move of some Mixed languages
28 * fix to phoenician-punic+ugaritic
29 * New hh.bib + clf ref updates
30 * synched refs
```

## 31 6. Conclusion

32 We have described the practices and principles for glottocodes as the identificational system for the languages, dialects and families of the world including data curation, technical infrastructure and update/version-tracking systematics. The resulting data observes the crucial aspects of the FAIR (Findable, Accessible, Interoperable, Reusable) Guiding Principles for scientific data management and stewardship. As such the glottocode-system responds to an important challenge in the realm of Linguistic Linked Data with numerous NLP applications.

## References

- 1 [1] H. Hammarström, R. Forkel, M. Haspelmath and S. Bank, Glottolog 4.4, 2021.
- 2 [2] M. Cysouw and J. Good, Languoid, Doculect, Glos-sonym: Formalizing the notion "language", *Language Documentation and Conservation* **7** (2013), 331–359.
- 3 [3] W. Schmidt, Gliederung der australischen Sprachen, *Anthropos* **7**, **7**, **8**, **9**, **12/13**, **12/13** (1912, 1912, 1912, 1913, 1914, 1917/1918, 1917/1918), 230–251, 463–497, 1014–1048, 526–554, 980–1018, 437–493, 747–817.
- 4 [4] Č. Loukotka, *Classification of the South American Indian Languages*, Reference Series, Vol. 7, Los Angeles: Latin American Center, University of California, 1968.
- 5 [5] G. van Bulck, *Les recherches linguistiques au Congo Belge: résultats acquis, nouvelles enquêtes à entreprendre*, Mémoires de l'Institut Royal Colonial Belge, Bruxelles, Vol. 16, Bruxelles: Institut Royal Colonial Belge, Bruxelles, 1948.
- 6 [6] J. Good and C. Hendryx-Parker, Modeling Contested Categorization in Linguistic Databases, in: *Proceedings of the EMELD 2006 Workshop on Digital Language Documentation: Tools and Standards: The State of the Art. Lansing, Michigan. June 20-22, 2006*, Lansing, Michigan: E-MELD, 2006, pp. 1–22.
- 7 [7] S. Nordhoff and H. Hammarström, Glottolog/Langdoc: Defining dialects, languages, and language families as collections of resources, in: *Proceedings of the First International Workshop on Linked Science 2011*, T. Kauppinen, L.C. Pouchard and C. Keßler, eds, CEUR Workshop Proceedings, Vol. 783, CEUR, 2011, pp. 1–7. <http://iswc2011.semanticweb.org/fileadmin/iswc/Papers/Workshops/LISC/nordhoff.pdf>.
- 8 [8] M. Mann and D. Dalby, *A thesaurus of African languages: a classified and annotated inventory of the spoken languages of Africa with an appendix on their written representation*, London: Hans Zell, London, 1987. ISBN 9780905450247.
- 9 [9] H. Hammarström, Ethnologue 16/17/18th editions: A comprehensive review, *Language* **91**(3) (2015), 723–737, Plus 188pp online appendix..
- 10 [10] D.M. Eberhard, G.F. Simons and C.D. Fennig, *Ethnologue: Languages of the World*, 24 edn, Dallas: SIL International, 2021. <http://www.ethnologue.com>.
- 11 [11] R. Forkel, Glottolog 3.0 released, 2017. <https://cldd.org/2017/03/29/glottolog-3-0.html>.
- 12 [12] R. Forkel, Glottolog 3.0 – A collaborative, versioned catalog of languages and dialects, 2016. <https://cldd.org/docs/poznan/glottolog-3-0.pdf>.
- 13 [13] R. Forkel, S.J. Greenhill and C. Rzymiski, glottolog/pyglottolog: Glottolog API, Zenodo, 2020. doi:10.5281/zenodo.3753876.
- 14 [14] M.D. Wilkinson, M. Dumontier, I.J. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J.W. Boiten, L.B. da Silva Santos, P.E. Bourne, J. Bouwman, A.J. Brookes, T. Clark, M. Crosas, I. Dillo, O. Dumon, S. Edmunds, C.T. Evelo,

- 1 R. Finkers, A. Gonzalez-Beltran, A.J. Gray, P. Groth,  
2 C. Goble, J.S. Grethe, J. Heringa, P.A. 't Hoen,  
3 R. Hooft, T. Kuhn, R. Kok, J. Kok, S.J. Lusher,  
4 M.E. Martone, A. Mons, A.L. Packer, B. Persson,  
5 P. Rocca-Serra, M. Roos, R. van Schaik, S.A. San-  
6 sone, E. Schultes, T. Sengstag, T. Slater, G. Strawn,  
7 M.A. Swertz, M. Thompson, J. van der Lei, E. van  
8 Mulligen, J. Velterop, A. Waagmeester, P. Wittenburg,  
9 K. Wolstencroft, J. Zhao and B. Mons, The FAIR  
10 Guiding Principles for scientific data management and  
11 stewardship, *Scientific Data* **3**(160018) (2016), 1–9.  
12 [15] J. Tension, G. Kellogg and I. Herman, Model for  
13 Tabular Data and Metadata on the Web, Technical  
14 Report, World Wide Web Consortium (W3C), 2015.  
15 <http://www.w3.org/TR/tabular-data-model/>.  
16 [16] R. Fielding and J. Reschke, Hypertext Transfer Proto-  
17 col (HTTP/1.1): Semantics and Content, RFC, 7231,  
18 RFC Editor, 2014. ISSN 2070-1721. [https://www.  
19 rfc-editor.org/rfc/rfc7231.txt](https://www.rfc-editor.org/rfc/rfc7231.txt).  
20 [17] F. Manola and E. Miller, RDF Primer, Technical  
21 Report, W3C, 2004. [https://www.w3.org/TR/2004/  
22 REC-rdf-primer-20040210/](https://www.w3.org/TR/2004/REC-rdf-primer-20040210/).  
23 [18] H. Hammarström, R. Forkel, M. Haspelmath  
24 and S. Bank, glottolog/glottolog-cldf: Glot-  
25 tolog database 4.3 as CLDF, Zenodo, 2020.  
26 doi:10.5281/zenodo.4061165.  
27 [19] R. Forkel, J.-M. List, S.J. Greenhill, C. Rzym-  
28 ski, S. Bank, M. Cysouw, H. Hammarström,  
29 M. Haspelmath, G.A. Kaiping and R.D. Gray,  
30 Cross-Linguistic Data Formats, advancing  
31 data sharing and re-use in comparative lin-  
32 guistics, *Scientific Data* **5**(180205) (2018), 1–  
33 10. doi:<https://doi.org/10.1038/sdata.2018.205>.  
34 <https://www.nature.com/articles/sdata2018205>.  
35 [20] R. Forkel, J.-M. List, M. Cysouw and S.J. Green-  
36 hill, CLDF 1.0, Technical Report, Max Planck Insti-  
37 tute for the Science of Human History, Jena, 2017.  
38 doi:10.5281/zenodo.1117644.  
39 [21] R. Pollock, J. Tension, G. Kellogg and I. Herman,  
40 Metadata Vocabulary for Tabular Data, Technical Re-  
41 port, World Wide Web Consortium (W3C), 2015.  
42 <https://www.w3.org/TR/tabular-metadata/>.  
43 [22] J. Tandy, I. Herman and G. Kellogg, Generating RDF  
44 from Tabular Data on the Web, Technical Report,  
45 World Wide Web Consortium (W3C), 2015. <https://www.w3.org/TR/csv2rdf/>.  
46 [23] M. Sporny, D. Longley, G. Kellogg, M. Lanthaler, P.-  
47 A. Champin and N. Lindström, JSON-LD 1.1 – A  
48 JSON-based Serialization for Linked Data, Technical  
49 Report, World Wide Web Consortium (W3C), 2020.  
50 <https://www.w3.org/TR/json-ld/>.  
51 [24] C. Lagoze, H.V. e Sompel, M. Nelson and S. Warner,  
The Open Archives Initiative Protocol for Meta-  
data Harvesting, Technical Report, Open Archives  
Initiative, 2002. [http://www.openarchives.org/OAI/  
openarchivesprotocol.html](http://www.openarchives.org/OAI/openarchivesprotocol.html).  
[25] E.M. Bender, On achieving and evaluating language  
independence in NLP, *Linguistic Issues in Language  
Technology* **6** (2011), 1–26.  
[26] D. Hovy and S.L. Spruit, The Social Impact of Nat-  
ural Language Processing, in: *Proceedings of the 54th  
Annual Meeting of the Association for Computa-  
tional Linguistics (Volume 2: Short Papers)*, Asso-  
ciation for Computational Linguistics, Berlin, Ger-  
many, 2016, pp. 591–598. doi:10.18653/v1/P16-2096.  
<https://aclanthology.org/P16-2096>.  
[27] H. Hammarström, R. Forkel, M. Haspelmath and  
S. Bank, glottolog/glottolog: Glottolog database 4.3,  
Zenodo, 2020. doi:10.5281/zenodo.4061162.