



CLLD – Cross-Linguistic Linked Data

Robert Forkel
Department of Linguistics, Max Planck Institute for Evolutionary Anthropology

Goal

Help collecting the world's language diversity heritage by providing interoperable data publication structures

Linguistic databases on the web

Observations

- ✓ Almost all quantitative papers at ALT 10 used WALS data.
- ✓ Many typologists at ALT 10 have their own typological database.

CLLD – The strategy

Since reuse tends to be the determining factor that keeps resources from vanishing, we want to bridge the gap between data collection and data reuse by

- ✓ publishing databases thereby incentivizing researchers through recognition;
- ✓ using technology that maximizes exposure of our data in the emerging web of data.

CLLD – The implementation

This twofold strategy is implemented by three service components:

- ✓ infrastructural: Glottolog - a comprehensive language catalog and bibliography,
- ✓ structural: Dictionaria - a dictionary journal and CrossGram Journal - a journal publishing typological databases,
- ✓ technological: c11d - a software platform for implementing linguistic database applications, which will be used to serve standalone database publications like IDS, WOLD, ASJP, WALS, APiCS, eWAVE, ValPal as well as the journals.

To maximize resuability

- ✓ we provide the data under Open Data Licenses,
- ✓ and the platform as Open Source software under a free license.

Linked Data - disseminating data in standard formats

- ✓ Defines a unified data access protocol for the web.
- ✓ Well-suited for distributed data providers
 - identifiers are URLs which are globally unique,
 - RDF and OWL provide the vocabulary to merge resources.
- ✓ Provides an easy to implement lowest level of service in a graceful degradation scenario

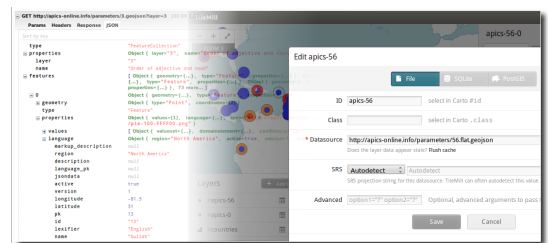
Use off-the-shelf tools to explore a dataset

Linked Data Explorer accessing Glottolog Linked Data serialized as RDF/XML.



Use off-the-shelf tools to transform a dataset

The map-making software Tilemill accessing APiCS data in GeoJSON format.



The data model

This implementation plan is realistic, because the underlying data model for the target databases is both reasonably simple and abstract enough to cover many use cases.

The data model comprises the following entities: Dataset, Contribution, Contributor, Language, Parameter, Value, UnitParameter, Unit, Sentence, Source.

Value – an APiCS datapoint

Datapoint Kriol/Pronoun conjunction

Both the juxtaposed inclusory construction and the conjoined construction are attested in Kriol. Both are relatively rare so the judgment of relative importance here is very impressionistic.

Values

- Inclusory pronoun juxtaposed with subset NP

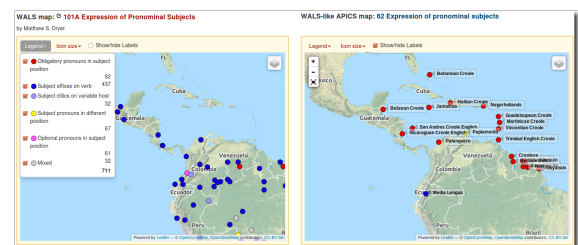
Show/hide details

Example 25-96 Sentence

Mindubala Namij kolim dardaga.
Mindubala Namij kol-im dardaga.
10U.KOZ. Namij kol 10 plant species
'Namij and I (in our language, Ngarinyman) call it dardaga (an edible plant).'

Contributor
Kriol - Eve Schultze-Bernold and
Denise Angelo cite
Feature
Pronoun conjunction
Source
Fieldwork Schulze-Bernold
Angelo et al. 1988a Source

Comparing WALS and APiCS data pulled in as GeoJSON



Unit – a Dictionaria word

Unit caa

deer
00:00 / 00:00

the deer (Hoosack Dictionary Dictionary)

Values

semantic domain
animal_mammal
semantic domain
physical_consumption_food
part of speech
noun

Source Contribution

Hoosack Dictionary by Iren
Hartmann cite



Mission accomplished? – A week's Visitors to Glottolog

