

Concepticon: A Resource for the Linking of Concept Lists

Johann-Mattis List¹, Michael Cysouw², Robert Forkel³

¹CRLAO/EHESS and AIRE/UPMC, Paris, ²Forschungszentrum Deutscher Sprachatlas, Marburg,

³Max Planck Institute for the Science of Human History, Jena
mattis.list@lingpy.org, emailpending, emailpending

Abstract

We present an attempt to link the large amount of different concept lists (aka “Swadesh lists”) which are used in the linguistic literature. This resource, the *Concepticon* (<http://concepticon.clld.org>), links **xxx** concepts labels from **xxx** conceptlists to **xxx** concept sets. Each concept set is given a unique numerical identifier, a unique label, and a human-readable definition. Concept sets are further structured by defining different relations between the concepts. The resource can be used for various purposes. Serving as a rich reference for new and existing databases in diachronic and synchronic linguistics, it allows researchers a quick access to studies on semantic change, cross-linguistic polysemies, and semantic associations.

Keywords: concepts, concept lists, Swadesh lists, cross-linguistically linked data

1. Introduction

In 1950, Morris Swadesh (1909 – 1967) proposed the idea that certain parts of the lexicon of human languages are universal, stable over time, and rather resistant to borrowing. As a result, he claimed that these parts of the lexicon, which was later called *basic vocabulary*, would be very useful to address the problem of subgrouping in historical linguistics:

[...] it is a well known fact that certain types of morphemes are relatively stable. Pronouns and numerals, for example, are occasionally replaced either by other forms from the same language or by borrowed elements, but such replacement is rare. The same is more or less true of other everyday expressions connected with concepts and experiences common to all human groups or to the groups living in a given part of the world during a given epoch. (Swadesh, 1950, 157)

He illustrated this by proposing a first *list of basic concepts*, which was, in fact, nothing else than a collection of concept labels, as shown below:¹

I, thou, he, we, ye, one, two, three, four, five, six, seven, eight, nine, ten, hundred, all, animal, ashes, back, bad, bark, belly, big, black, blood, bone, brother (elder), child (son or daughter), cloud, cold, come, cry (weep), dance, day, dog, dust, ear, earth, eat, egg, eye, far, father, fire, flower, fog, foot, good, grass, green, hair, hand, head, heart, here, hit (with fist), hunt, husband, ice, lake, laugh, leaf, left hand, leg, liver, long, louse, man, meat, mother, mountain, mouth, name, near, neck, night, nose, person, rain, red, right

hand, road (trail), root, rope, salt, sand, short, sing, sister (elder), skin, sky, small, smoke, snake, snow, speak, spear (war), star, stone, sun, swim, tail, that, there, this, tongue, tooth, tree, warm, water, what, where, white, who, wife, wind, woman, year, yellow. (Swadesh, 1950, 161)

In the following years, Swadesh refined his original concept lists of basic vocabulary items, thereby reducing the original test list of 215 items first to 200 (Swadesh, 1952) and then to 100 items (Swadesh, 1955). Scholars working on different language families and different datasets provided further modifications, be it that the concepts which Swadesh had proposed were lacking proper translational equivalents in the languages they were working on, or that they turned out to be not as stable and universal as Swadesh had claimed (Matisoff, 1978; Alpher and Nash, 1999). Up to today, dozens of different concept list have been compiled for various purposes. They are used as heuristical tools for the detection of deep genetic relationships among languages (Dolgopolsky, 1964), as basic values for traditional lexicostatistical and glottochronological studies (Dyen et al., 1992; Starostin, 1991), or as litmus test for dubious cases of language relationship which might be due to inheritance or borrowing (McMahon et al., 2005; Chén Bǎoyà 陈保亚, 1996; Wang, 2006).

Apart from concept lists proposed for the application in historical linguistics, there is furthermore a large amount of not explicitly diachronic data, including concept lists serving as the basis for field work in specific linguistic areas (Kraft, 1981), concept lists which serve as the basis for large surveys on specific linguistic phenomena (Haspelmath and Tadmor, 2009), or concept lists which deal with the internal *structuring* of concepts, be it cognitive associations (Nelson et al., 2004), cross-linguistic polysemies (Lis, 2014), or frequently recurring semantic shifts (Bulakh et al., 2013).

¹This list contains 123 items in total. According to Swadesh, these items occurred both in his original test list of English items, and in the data on the Salishan languages, which he employed for his first glottochronological study.

2. Concept Lists

Concept lists are – simply speaking – lists of concepts. In these lists, concepts are ideally described with help of a *concept label* and also a short *definition*. Most published concept lists, however, only contain a concept label. On the other hand, certain concept lists have been further expanded by adding structure, such as *rankings*, *divisions*, or *relations*.

2.1. Purpose of Concept Lists

Concept lists are compiled for a variety of different purposes. A major distinction can be made between those concept lists which have been compiled for the purpose of *language comparison* and those which have been compiled for the purpose of *concept comparison*. Among the former, we can further distinguish those lists which are used to prove *language relationship* (Dolgopolsky, 1964), those which are used for *linguistic subgrouping* (Norman, 2003; Starostin, 1991; Swadesh, 1955), and those which can be used to identify *contact layers* (Chén Bǎoyà 陈保亚, 1996). Among the latter, we can distinguish between concept lists with a primarily *synchronic objective* (Hill et al., 2014), and those with a primarily *diachronic objective* (Haspelmath and Tadmor, 2009; Bulakh et al., 2013).

2.2. Structure of Concept Lists

The purpose for which a given concept list was originally defined has an intermediate influence on its structure. Given the multitude of use cases in both synchronic and diachronic linguistics, it is difficult to give an exhaustive and unique classification schema for all concept lists which have been compiled in the past. In Table 1, we have nevertheless tried to distinguish eight basic types of concept lists and give one list for each of the types as a prototypical example.²

Type	Example	Purpose
basic vocabulary list ("Swadesh list")	Swadesh 1952 / 200 items	subgrouping
subdivided concept list	Yakhontov 1991 (see Starostin 1991) / 35 + 65 items	genetic relationship, layer identification
"ultra-stable" concept list	Dolgopolsky 1964 / 15 items	genetic relationship
questionnaire	Allen 2007 / 500 items	dialect / language comparison
ranked list	Starostin 2007 / 110 items	subgrouping, layer identification
list of concept relations	DatSemShift, Bulakh et al. 2013 / 2424 items	representation of concept relations
special-purpose concept list	Matisoff 1978 / 200 items	subgrouping of Tibeto-Burman languages
historical concept list	Leibniz 1768 / 128 items	language comparison

Table 1: Examples for different types of concept list as they can be found in the literature.

3. Linking Concept Lists

Here, we should describe the basic characteristics of the concepticon, like the way the things are linked with each other. Maybe, including a graphic would

²For further information regarding these concept lists, just click on the links in the "Example" field of the table.

also be useful The concepticon is an attempt to link the many different concept lists ("Swadesh Lists") which are used in the linguistic literature. In practice, all entries from the various concept lists are linked to a *concept set* as an intermediate way to reference the concepts. The Concepticon currently links **xxx** concepts from **xxx** concept lists to **xxx** concept sets and defines **xxx** relations between the concept sets.

A concept list is a collection of concepts that is deemed interesting by scholars. Minimally, it consists of an *identifier* for each concept which the lists contains, and a *label* by which the concept is referenced. The creator of a concept list is called a *compiler*. Each concept list is tight to one or more *sources*, it is given in one or more *source languages* and was compiled for one or more *target languages*. A *description* gives further information on each concept list in free, exclusively human-readable form.

To facilitate our workflow and to guarantee the comparability of concept lists even if they do not share concepts which are directly linked via our concept sets, we define additional and very simple *concept relations* between concept sets (*broadier*, *narrower*, *similar*). Even if the concepts in two or more concept lists are not assigned to the same concept set, they can still be assigned to concept sets via concept relations.

4. Examples

Examples may be useful to be included, they could, however, also be put into a nice graphic in which the web-application is presented along with the underlying graphs showing the relations between the concepts

4.1. Rain

Maybe give rain as an example here, as in the slides

4.2. Child

Child example, for hierarchies

4.3. Burn

Burn example to show problems with transitivity etc.

5. Using the Concepticon

Hier eventuell zeigen, wie das Concepticon bei Dictionaria und Lexibank benutzt werden kann

6. Outlook

hier noch mal sagen, dass wir natürlich noch weiter daran arbeiten.

7. Acknowledgements

vielleicht nicht nötig im Moment...

8. References

- Allen, B. (2007). *Bai Dialect Survey*. SIL International.
- Alpher, Barry and Nash, David. (1999). Lexical replacement and cognate equilibrium in australia. *Australian Journal of Linguistics: Journal of the Australian Linguistic Society*, 19(1):5--56.
- Bulakh, M., Ganenkov, Dimitrij, Gruntov, Ilya, Maisak, T., Rousseau, Maxim, and Zalizniak, A. (2013). Database of semantic shifts in the languages of the world.
- Chén Bǎoyà 陈保亚. (1996). *Lùn yǔyán jiēchù yǔ yǔyán liánméng*. Yǔwén 语文, Běijīng 北京.
- Dolgopolsky, Aron B. (1964). Gipoteza drevnejšego rodstva jazykovych semej Severnoj Evrazii s verojatnostej točki zrenija [A probabilistic hypothesis concerning the oldest relationships among the language families of Northern Eurasia]. *Voprosy Jazykoznanija*, 2:53--63.
- Dyen, Isidore, Kruskal, Joseph B., and Black, Paul. (1992). An indoeuropean classification. *Transactions of the American Philosophical Society*, 82(5):iii--132.
- Haspelmath, Martin and Tadmor, Uri. (2009). World loanword database.
- Hill, Felix, Reichart, Roi, and Korhonen, Anna. (2014). Simlex-999: Evaluating semantic models with (genuine) similarity estimation. *CoRR*, abs/1408.3456.
- Kraft, Charles H., editor. (1981). *Chadic wordlists*. Dietrich Reimer, Berlin.
- (2014). Clics: Database of Cross-Linguistic Colexifications.
- Matisoff, James A. (1978). *Variational semantics in Tibeto-Burman. The 'organic' approach to linguistic comparison*. Institute for the Study of Human Issues.
- McMahon, April, Heggarty, Paul, McMahon, Robert, and Slaska, Natalia. (2005). Swadesh sublists and the benefits of borrowing: An Andean case study. *Transactions of the Philological Society*, 103:147--170.
- Nelson, Douglas L., McEvoy, Cathy, and Schreiber, Thomas A. (2004). The university of south florida free association, rhyme, and word fragment norms. *Behaviour Research Methods, Instruments, & Computers*, 36(4):402--407.
- Norman, Jerry. (2003). The Chinese dialects. In Thurgood, G. and LaPolla, R., editors, *The Sino-Tibetan languages*, pages 72--83. Routledge, London and New York.
- Starostin, Sergej Anatol'evic. (1991). *Altajskaja problema i proischoždenije japonskogo jazyka [The Altaic problem and the origin of the Japanese language]*. Nauka, Moscow.
- Swadesh, Morris. (1950). Salish internal relationships. *International Journal of American Linguistics*, 16(4):157--167.
- Swadesh, Morris. (1952). Lexico-statistic dating of prehistoric ethnic contacts. *Proceedings of the American Philosophical Society*, 96(4):452--463.
- Swadesh, Morris. (1955). Towards greater accuracy in lexicostatistic dating. *International Journal of American Linguistics*, 21(2):121--137.
- von Leibniz, Gottfried Wilhelm. (1768). Desiderata circa linguas populorum, ad dn. podesta. In Dutens, Louis, editor, *Godefridi Guilielmi Leibnitii opera omnia, nunc primum collecta, in classes distributa, praefationibus et indicibus exornata*, volume 6.2, pages 228--231. Fratres des Tournes, Geneva.
- Wang, Feng. (2006). *Comparison of languages in contact. The distillation method and the case of Bai*. Institute of Linguistics Academia Sinica, Taipei.