

A uniform data model for joint publication of minor language dictionaries

Martin Haspelmath
2014

1. Introduction

This paper proposes a data model for dictionaries of minor languages that can be applied to any language and that allows publication of many different dictionaries in the same general framework, with identical search options and linking of entries via comparison meanings. Dictionaries of minor languages (i.e. languages which are not taught widely in schools and do not have official support) are of great interest to linguistic scholarship, but publishing such dictionaries in a way that that permits easy use by linguists and speakers requires electronic publication in database form and a relatively uniform data structure.¹

The basic idea is that it should not be necessary to create a separate digital infrastructure for each new electronic dictionary. While it is possible to create separate infrastructures for commercial dictionaries of big languages such as merriam-webster.com, dudende.de and lrousse.com, this is not a realistic option for languages that have just a few thousand speakers or less and a handful of linguists working on them. Thus, electronic dictionaries of minor languages must adopt the Wikipedia model of having a single infrastructure for joint publication (e.g. a dictionary journal), and this implies a UNIFORM DATA MODEL.

By *data model* we mean a set of data types and their interrelations, abstracting away from the way in which they are presented visually. Our data model should be able to subsume all the important information found in traditional linear dictionaries, but since it must be uniform and should be relatively simple, the information will not always have the same form as traditional dictionaries. In particular, information that is conveyed by purely linear arrangement must be represented in some different form in our model, because databases do not have a linear structure.

It is immediately apparent that all dictionaries share some key properties: they consist of lists of WORD FORMS, each of which is linked to some additional information about these word forms, in particular a MEANING DESCRIPTION in a metalanguage, a meaning description in the object language itself, GRAMMATICAL INFORMATION such as word-class, gender or inflection class, and so on. Dictionaries also usually contain EXAMPLE SENTENCES that illustrate the use of a word in context. Thus, it should be possible in principle to adopt the same data model for all languages, regardless of language-specific peculiarities of individual languages.

This paper describes such a data model and proposes specific solutions for data modeling issues that dictionary authors may be confronted with. Of course, by using a

¹ See Bell & Bird (2000) for a similar attempt to devise a uniform data model. However, their focus is on formal implementation (in XML) and on capturing everything that existing linear dictionaries can be found to contain. As a result, their model is far more complicated than the model proposed here. The focus here is on a data model for future dictionaries that are created with a view to integrating them into a larger set of electronic dictionaries.

single data model for all languages, we may sometimes have to adopt a solution that is less than perfect for a particular language. But we think that using different data models will make it far more difficult to store and access the information, so that it is best to accept the limitations of the uniform data model.

2. Data types

Dictionaries are collections of LEXICAL ENTRIES. The most important parts of an entry are the LEMMA and the SENSE. The lemma is the citation form which serves as the name of an entry. In the simplest form, a dictionary is a list of lemma-sense pairs (*arbre*: ‚tree‘, *maison*: ‚house‘, *acheter*: ‚buy‘, etc.).

The most typical lexical entries are SIMPLE LEXEMES such as nouns and verbs, but more complex lexical entries, such as PHRASAL ENTRIES (*take part*, *point of view*, etc.) and IDIOMS (*all hell broke loose*, *fall on deaf ears*) must also be treated as lexical entries.

A lexical entry may be associated with one or more senses. While the simplest case of words with a single sense (e.g. *estomac*: ‚stomach‘) is not uncommon, there are also many cases of words with multiple senses. Thus, the data model needs to accommodate one-to-many correspondences between lemmas and senses, e.g. German *Platz*:

Platz (i) ‚space‘
 (ii) ‚seat‘
 (iii) ‚square‘

Both lexical entries (§4) and senses (§5) are associated with further categories of information (database fields).

3. Database fields

For concreteness, the categories of information that are attached to lexical entries and senses are called DATABASE FIELDS. The most important question that a dictionary author needs to resolve is which database fields to use. In traditional print dictionaries, the various categories of information are printed in a linear fashion, with no labels for the database fields, and mostly only with typographic distinctions or separation marks such as commas, semicolons or brackets between the fields.

For example, Figure 1 shows an online dictionary of Maung (Australian), where lexical entries have the database fields WORD-CLASS (blue italics, abbreviated) and RELATED ENTRY (green, preceded by *See*:), and each entry can have multiple senses. A sense is a list of one or more SENSE DESCRIPTORS (e.g. the two descriptors „gather together in one place“ and „put together“ for the single sense of the entry *akpaj*). Each sense may be associated with a SEMANTIC CATEGORY (purple, preceded by *Category*:) as well as with an EXAMPLE (blue).

ajput <i>nc4</i> .
1 • sand, beach. <i>Category: Elements. Ngana tuka ajput.</i> I am going to the beach.
2 • sugar. <i>Ngana nganalakpalwarrki ngartu ta ajput.</i> I am going to buy myself sugar.
ajputajput <i>nc4</i> . along the beaches. <i>Kawarra ajputajput.</i> They are going along the beaches (looking for turtle eggs).
akarnpa <i>nc4</i> . yam root. <i>Category: Food Drink Cooking and Fire. Takapa akarnpa wularr.</i> Those are yam roots. <i>See: igarnpa.</i>
akarra <i>nc6</i> . many, a lot of. <i>Category: Number and Quantity. See: wigarra.</i>
akijalk <i>nc6</i> . yam, grain or vegetable fruit or body. <i>Category: Plants. See: ingijalk.</i>
akiri <i>nc6</i> . skin of some <i>nc6</i> item, yam, grain, etc., foreskin. <i>Category: Plants. See: ingiri.</i>
akpaj <i>vroot</i> . gather together in one place, put together. <i>Gram: used with tv kinima Category: Holding and Transfer. Ararrkpi ying inimang akpaj.</i> A man might come and put them together.

Figure 1: AuSIL Dictionary of Maung²

Thus, for the dictionary in Figure 1, the following database fields are used:

lexical entry:

- lemma (= citation form) (black, boldface)
- word-class (= part of speech) (blue, italics, serif font)
- associated entry (green)

sense:

- (– sense number)
- list of sense descriptors (black serif font, separated by comma)
- semantic category (purple)
- example (dark blue, translation in black serif font)

The list and organization of the database fields that are used by a dictionary is its DATA MODEL. Different dictionaries tend to have slightly different data models, but here a single uniform data model is proposed that should be able to accommodate most of the information that authors want to provide in dictionaries of minor languages.

Some fields are applicable to any language in principle (GENERAL FIELDS), while other fields may be language-specific (SPECIFIC FIELDS). The fields used in the Maung dictionary in Figure 1 are all general fields, i.e. they could be used also in a dictionary of any other language. But dictionaries may also contain specific fields that are useful only for the language in question, e.g. information on inflectional class membership.

² <http://ausil.org/Dictionary/Maung/lexicon/index.htm> (a dictionary created with the programme Lexique Pro)

4. Database fields for lexical entries

After laying out the main concepts, the following sections give the uniform data model that is proposed here.

4.1. General fields

The following database fields of lexical entries are immediately clear and do not raise any questions. For dictionaries of minor languages, lemma and part-of-speech are generally considered essential and are present in all dictionaries (hence in boldface below).

- **lemma = headword = citation form**³
- lemma in original script (if the language is commonly written in another script)
- pronunciation of lemma (in IPA)
- **part-of-speech**
- general comments
- bibliographical reference

In addition, each entry is obligatorily associated with one or more senses (see §5 for data fields of senses).

Lexical entries may also contain historical information, especially on loanwords, but also other etymological information:

- source language (for loanwords)
- source word (for loanwords) plus gloss
- etymology

It should be noted that lexical entries need not be simple words, but can be idioms of various sizes (e.g. *hot dog*, *pay lipservice*, *once and for all*, *be that as it may*). Such complex entries can be associated with their parts via association fields (§6).

4.2. Specific fields

In addition to the general fields, dictionaries of particular languages may require other kinds of fields, in particular fields for grammatical information:

- verbal valency
- gender
- inflectional class⁴

³ Several different entries may have the same lemma (such cases are called HOMONYMS). For easier distinguishability, they may be accompanied by a homonym number that is part of the lemma (e.g. *bank 1*, *bank 2*)

⁴ Inflectional information is often given in the form of specific inflected forms, e.g. unpredictable plural or past-tense forms (called „principle parts“ if these can be used to predict further inflected forms). Such specific forms are association fields, see §6 below.

Other specific fields may concern sociolinguistic information, e.g. concerning dialects, registers, politeness, taboo, speech level (plain vs. honorific), about literary vs. colloquial usage.

These fields cannot be specified in advance, because they may depend on language-specific circumstances. There is thus no way around language-specific fields for lexical entries.

5. Database fields for senses

The most important field for senses is of course the sense description (= definition). As mentioned earlier, a lexical entry can have several different senses, which are often numbered (e.g. in the Maung entry *ajput* in Figure 1).⁵

A sense can be a single sense descriptor, or can consist of a list of sense descriptors, when it is described by multiple words or phrases (e.g. the first sense of *ajput* is characterized by two sense descriptors as ‚sand, beach‘).

There is no hard and fast rule for when to set up a new sense, as opposed to describing a single sense with multiple descriptors. In general, multiple senses are used to represent POLYSEMY, i.e. a situation where a lexical form has a number of clearly distinct senses which cannot be readily subsumed under a single meaning from which the senses can be contextually derived. Thus, while both squares and seats are kinds of spaces, it is not possible to predict that German *Platz* ‚space, seat, square‘ (mentioned above in §2) should have precisely these three senses (and not, for instance, the sense ‚apartment‘, even though an apartment is also an important kind of space, cf. English *place*, which can have this sense). Deciding between these two kinds of representations is not straightforward at all in practice, but what is important for our data model is that both ways of representing a word’s meaning are often used by lexicographers and hence need to be recognized by the data model.

In addition to the sense description in the metalanguage (English), senses can be given in the object language itself (as in monolingual dictionaries), and in some other metalanguage that is relevant for the minor language in question (most likely a national language such as Spanish or Russian). Many dictionaries also give a semantic domain.

- **sense description (English)**
- sense description in object language („native definition“)
- sense description in another relevant metalanguage (e.g. Spanish, Russian)
- semantic domain
- comment

Finally but importantly, each sense may be associated with an EXAMPLE. Examples are particularly helpful when a word has multiple senses, and dictionaries tend to give

⁵ There is no hard and fast rule for when to set up homonyms and when to assume a single entry with multiple senses. In general, homonyms are set up when the grammatical information is different and/or when the senses are not related at all. Given our data model, when the grammatical information is different, it is necessary to set up a new entry (see §7).

examples especially for such polysemous words. An example may be a clause or a phrase (e.g. a verb with an object, or a complex noun phrase).

Examples minimally consist of the primary text and the free translation, but ideally they should also have an interlinear morpheme-by-morpheme translation.⁶ There may be multiple examples for each sense, and an example may illustrate multiple lexical entries, so there is a many-to-many relation between examples and senses. Examples thus need to be in a separate database table.

It appears that the data fields for senses are all general, i.e. there are no language-specific fields here.

6. Association fields

One important type of information for many lexical entries (and senses) is association with some other lexical entry. For example, we may want to say that a lexical entry is part of, or contains, or is derived from, some other lexical entry. For example:

(lexical entry:) *cloth*
 – part of (lexical entry): *table cloth*

(lexical entry:) *tablecloth*
 – contains (lexical entries): *table, cloth*

(lexical entry:) *hot dog*
 – contains (lexical entries): *hot, dog*

(lexical entry:) *amusement*
 – derived from (lexical entry): *amuse*

For senses, we may say that they are synonymous with other senses, or antonyms, or hyponyms, etc.

(sense:) ‚stop’ (a sense of the entry *cease*)
 – synonym(ous with): ‚stop’ (a sense of the entry *stop*)

Association fields are also used to represent unpredictable inflectional information, or information about associated classifiers etc.

(lexical entry:) *mouse*
 – plural form: *mice*

⁶ It is mainly for typographic reasons that examples in dictionaries of minor languages tend not to have interlinear translations, whereas this is now completely normal for grammars of minor languages. But they are of course no less necessary in dictionaries, because readers cannot be assumed to know the language well enough to understand the example without such interlinear glosses.

(lexical entry:) *mǎ* ‚horse’ (Mandarin Chinese)
 – numeral classifier: *pǐ*

Some associated entries such as inflected forms like *mice* need not be given a full lexical entry. *Mice* would be an ASSOCIATED ENTRY, which does not have word-class and sense information, but only has the association field „plural form of“, which links it to the regular entry *mouse*.

In electronic dictionaries with hypertext links, such associations can simply be implemented by links to the corresponding lexical entries (and senses) via the lemmas.

7. One-to-many relationships

As was noted in §2, a lexical entry may be associated with multiple senses, as is again illustrated below.

lexical entry: *macinare* (Italian)
 sense: grind
 sense: spend (a lot of money)
 sense: eat (a lot of food)

And of course a sense can be associated with multiple examples. Moreover, examples may of course be associated with multiple senses of different entries, illustrating simultaneously the uses of several words.

Multiple senses of an entry must occur in some linear order, and this linear order is sometimes given some significance in traditional dictionaries. For instance, Hausa *dari* has the two senses ‚1. dry cold’ and ‚2. chills due to illness’ (Newman’s Hausa dictionary, see Figure 3 in §9.x below). Presumably the ‚dry cold’ meaning is given first here because the other meaning represents some kind of semantic extension. In order to capture the significance of linear order, our data model would have to add a rank field to each sense. It is not clear whether this is necessary, because in many cases, senses seem to be linearized (and numbered) without implying any rank (primary vs. secondary meaning).

The association relations described in §6 may also be one-to-many. For example, a lexical entry can be said to comprise several other lexical entries (as in the case of *tablecloth*, which comprises *table* and *cloth*).

However, other one-to-many (or many-to-many) relationships need not be recognized. For example, it is not necessary to allow entries with two different parts of speech, such as:

lexical entry: *stop* (English)
 part-of-speech: verb
 part-of-speech: noun

lexical entry: *lingüista* (Spanish)
 gender: masculine
 gender: feminine

Linguists might want to posit precisely such one-to-many relationships, but such entries would lead to too much complexity in the data model, which is rarely needed for dictionaries of minor languages. If there are two different parts-of-speech, then these two are presumably also associated with two different senses, so that one would have to express that the lemma has sense 1 if it has part-of-speech 1, and sense 2 if it has part-of-speech 2, and so on. This would effectively mean the creation of a subentry, but we do not allow subentries, as discussed in §8. In the case of nouns with two genders, an alternative (if one wants to avoid multiple lexical entries) might be to set up a new gender category „epicene“.

8. No subentries

Our data model does not allow for subentries. Traditional dictionaries often have a structure such as:

put
 put aside
 put off
 put up

That is, the formally complex entries *put aside*, *put off* and *put up* are subentries of the formally simple entry *put*. This makes sense in linearized, paper-based dictionaries, where searching is done by alphabetic scanning.

In electronic dictionaries, this does not make much sense, because complex entries otherwise have the same properties as simple entries. The important piece of information that is traditionally expressed by the relationship between main entry and subentry is the part-of and the contains relationship:

entry: *put*
 – part of: *put aside*

entry: *put aside*
 – contains: *put, aside*

9. Some example dictionaries and how our data model applies to them

9.1. Kari’s Degema dictionary

In Ethelbert Kari’s (2008) dictionary of Degema, each entry is followed by the pronunciation in IPA and the word-class, plus one or more numbered senses (three senses occur with *kpor*, for example), and an example.

kpomoy [kpomój] <i>vi</i> be ill many times/always <from kpom >	kpotumañinesey [kpotuməñineséj -k-] <i>vt</i> cause to be crowded always <from kpotumañine >
kpomony DT [kpomónj] <i>vt</i> see kuroy	kpow [kpów] <i>vt</i> prevent (esp. the edges of a basket, cloth, etc.) from getting loose
kpooow [kpoóow] <i>interj</i> a response to a delayed answer to a call when the person called eventually answers	kpɔ¹ [kpɔ́] <i>vi</i> be old, be mature: Igbény yɔ ɪkpóte 'The mangoes are mature'.
kpɔr [kpór] <i>vt</i> 1 beat (esp. a drum) 2 play (esp. music) 3. sing	kpɔse [kpɔsé] <i>vt</i> cause to be old/mature <from kpɔ¹ >
kpɔrone [kpɔroné] <i>vt</i> sing about oneself <from kpɔr >	kpɔsey [kpɔséj] <i>vt</i> cause to be old/mature always <from kpɔ¹ >
kpɔñine [kpɔñiné -k-] <i>vt</i> sing about each other <from kpɔr >	kpɔviriy [kpɔβiríj -v- -t-] <i>vi</i> be old/be mature always <from kpɔ¹ >
kpɔroy [kp'rój] <i>vt</i> beat (esp. a drum)/play (esp. music)/sing many times/always <from kpɔr >	kpɔ² [kpɔ́] <i>vi</i> (of a bat, etc.) hang down from a tree/tree branch: Igbom ɪkpó m'ɪnwíny útány yɔ 'Bats are hanging down from the tree'.
kpɔron¹ [kpɔrón] <i>vt</i> remove the carapace of a crab/cover of a book, etc.	kpɔviriy [kpɔβiríj -v- -t-] <i>vi</i> (of a bat, etc.) hang down from a tree/tree branch always <from kpɔ² >
kpɔrone [kpɔroné] <i>vi</i> (of the carapace of a crab/cover of a book, etc.) be removed <from kpɔron¹ >	

Figure 2: An excerpt from Kari (2008) (A dictionary of Degema)

What is notable in Figure 2 is that many entries have subentries for morphologically related items, e.g. *kpɔ 1* (‘be old’), which has the subentries *kpɔse* (‘cause to be old’), *kpɔsey* ‘cause to be old always’, and *kpɔviriy* ‘be old always’. Since our model does not allow subentries (§8), these would have to be separate entries, but they would be linked to the entry *kpɔ 1* by association fields (§6), e.g.

(lexical entry:) *kpɔsey* ,cause to be old always’
 – contains root: *kpɔ 1* ,be old’

(lexical entry:) *kpɔ 1* ,be old’
 – habitual form: *kpɔviriy* ,be old always’

9.2. Newman’s Hausa dictionary

In Paul Newman’s (2007) dictionary of Hausa (see Figure 3), most entries have a word-class (in italics: *m*, *f*, *id*, *num*, *v1*, *v2*, etc), and quite a few have different senses, especially the verbs. Examples of multiple sense descriptors for a single sense are ‘honor, glory’ (one sense of the entry *dàukaka*), and ‘lift up, carry, take away’ (one sense of *dauka*). Quite a few of the entries have example phrases or clauses. Some of the nouns have unpredictable plural forms given in angle brackets. Occasionally there are sociolinguistic annotations such as „(obs)“ (obsolete) for *dari*².

Since Hausa verbs typically come in „grades“ of semantically related verbs derived from the same root, all verbs are treated as subentries of verb roots in the dictionary. Thus, the verbs *daura* ‘tie’ (grade 1), *daure* ‘tie up’ (grade 4), and *dauro* ‘happen to be’ (grade 6) are all found under the entry *DAUR-*. In our datamodel, which does not allow subentries, the verb root would be an associated entry (like English *mice*, cf. §6), and it would be linked to the individual verbs (and vice versa) via association fields.

<p>dār-dār <i>id</i> Palpitation of the heart (due to fear or anxiety), apprehension: Gābānā yanà ta ~. I was very nervous, full of anxiety.</p> <p>dārē <i>m</i> 1. hawan ~ Horseback riding with both legs on one side. 2. zaman ~ Rigid sitting of a chief indicating that he is in an officious mood.</p> <p>dāri¹ <i>num</i> <darurrukā> Hundred.</p> <p>dāri² <i>m</i> <darurrukā> (obs.) Halfpenny in old Nigerian currency.</p> <p>dāri <i>m</i> 1. Dry cold (<i>cf.</i> sanyī). 2. Chills due to illness.</p> <p>dāri-dāri <i>m</i> (usu. followed by gen. linker plus neg. subjunctive or v.n.) Timidity, apprehensiveness, being concerned about: Inā ~n kadà yà màkarà. I am concerned that he will come late.</p> <p>dāriḥā <i>f</i> Mystical Islamic sect, specifically the Tijaniyya sect.</p> <p>DARS- darsà <i>vI</i> (usu. used in ~ i.o. à rāi) Occur to s.o. — darsu <i>v7</i> (with à zūciyā) Occur to s.o.: Yā ~ à zūciyā. It became clear to me.</p> <p>darurrukā <i>pl. of dāri.</i></p> <p>dātā <i>f</i> A small green bitter tomato.</p> <p>dātā-dātā <i>f</i> A bitter grass (= dātārniyā).</p> <p>dātātā <i>pl. of dātā-dātā.</i></p> <p>dau ⇒ DAUK-</p> <p>DAUK- duka <i>v2</i> 1. Take. 2. Lift up, carry, take away: Barcī vā dūkē tà. She fell</p>	<p>holiday.</p> <p>DAUKAK- dukaḥā <i>vI</i> 1. Lift up, raise up. 2. ~ kārā Make a legal appeal. 3. Honor, exalt, promote.</p> <p>daukaka <i>f</i> 1. Honor, glory. 2. Promotion.</p> <p>daukàn-idò <i>adj</i> Dazzling.</p> <p>daukà-wuyà <i>m</i> Carrying s.o. on the shoulders.</p> <p>daukē <i>m</i> Suppository, usu. for children.</p> <p>daukī <i>m</i> Distributing to members of a household their proper share of food.</p> <p>dauki-bā-dadī <i>m</i> Combat, confrontation, struggle back and forth.</p> <p>dauki-daidai <i>m</i> 1. Single elimination. 2. Stealing things one by one, esp. items such as peanuts, kolanuts, or mangoes that are displayed for sale in a pile.</p> <p>DAUR- daurà <i>vI</i> 1. Tie sth onto sth: Nās yā ~ minī bandējī à hannu. The nurse put a bandage on my arm. 2. ~ aurē Perform a marriage. — daurè <i>v4</i> 1. Tie up. 2. Arrest, imprison. 3. <i>Idiom</i>: ~ fuskà Scowl. 4. ~ i.o. gindī Support, back up s.o. — daurō <i>v6</i> Happen to be: Hakà tā ~ à yi. That is how it should be done.</p> <p>daurarrē <i>adj.pp</i> Tied, imprisoned. — <i>m</i> Prisoner, inmate.</p> <p>DAURAY- daurayē <i>v4</i> Rinse.</p> <p>daurayà <i>f</i> 1. Dishes that have been rinsed off. 2. <i>v.n. of daurayē.</i></p>
---	---

Figure 3: An excerpt from Newman's (2007) (A dictionary of Hausa)

A few further observations: (i) One entry (*dari-dari*) has a miscellaneous grammatical comment (in parentheses) that should probably be in a general comments field. (ii) The entry *dāre* does not have a sense, but is associated with two (what appears to be) fixed phrases, *hawan dare* and *zaman dare*. Since their meaning (apparently) cannot be compositionally derived, they are idioms and hence they should get their own entries, linked to the association entry *dāre*. (iii) The dictionary treats some homonyms as senses, e.g. *dauraya* is said to have two senses: 'dishes that have been rinsed off', and 'verbal noun of *daurayē*'. The latter should be an association entry in our datamodel, linked to the verb *daurayē* 'rinse', so it cannot be another sense of the same entry.

9.3. Nichols's Chechen dictionary

Johanna Nichols's Chechen dictionary gives all Chechen words both in the Cyrillic spelling and in the pronunciation (using a non-IPA transcription). Nouns have gender information in round brackets (*v:b, j:j, b:b* etc.), declension-class information in square brackets, as well as the principal parts ergative singular and nominative plural (e.g. *taam, toomuo, teemash*). Verbs have conjugation class information in square brackets and valency information in round brackets. In addition, principal parts are given for many verbs (present tense, witnessed past tense, e.g. *taqā, teqa, teqira*).

member, fellow clansman.	<i>luu, taaluush</i> . Hulk, large person, fat person.
тайпаныйиша <i>n.</i> (v:b) <i>taipanjisha</i> . Clan sister, fellow clanswoman.	талх <i>n.</i> (j:j) [1] талхаш. <i>talx, talxash</i> . Piece, hunk (of meat).
тайп-тайпана <i>adj.</i> <i>taip-taipana</i> . Different, diverse, various, variable.	талха <i>v.</i> (Nom) [VIII] телха, телхира. <i>talxa[~], telxa, telxira</i> . Spoil, rot.
тайша <i>adj.</i> <i>taisha</i> Brown (color of cows).	там <i>n.</i> (b:b) [3] томо, темаш. <i>taam, toomuo, teemash</i> . Blood price, wergeld; brideprice. ◦ то-мана <i>toomana</i> As a favor
така <i>n.</i> (d:d) [7] такано / тако, такнаш / стакнаш. <i>taka, taknuo / takuo, taknash</i> . Glass, drinking glass, tumbler.	тамашийна <i>adj.</i> <i>taamashiina</i> . Surprising, strange, odd, curious, peculiar.
такха <i>v. simul.</i> (Nom) [VIII] текха, текхира. <i>taqa[~], teqa, teqira</i> . Crawl, creep.	там бала <i>v.</i> <i>taam bala[~]</i> . Pay brideprice, pay blood price, pay wergeld.
такха <i>v. simul.</i> (Erg-Nom) [v] токху, текхира. <i>taqa[~], toqu, teqira</i> . Repay debt, compensate; recover (from illness).	там бан <i>v.</i> (Erg-Dat-# <i>taam</i>) <i>taam ba[~]</i> . Do a favor, pay a compliment, please; give a bribe.
такхор <i>n.</i> (b:d) [2] такхоро, такхораш. <i>taqor, taquoruo, taquorash</i> . Stack, rick (of cornstalks).	таммаг1а <i>n.</i> (d:d) [1] таммаг1наш. <i>taammagh, taammagh-nash</i> . Handprint, footprint, fin-

Figure 4: An excerpt from Nichols (2010) (A dictionary of Chechen)

Different senses are separated by semicolons (e.g. *taam* ,blood price, wergeld; brideprice’). Fixed phrases are given their own entries, as in our datamodel (e.g. *taam bala[~]* ,pay brideprice’).

The main challenge of this dictionary for our datamodel is that there are multiple cases where forms are given both in transcription and in Cyrillic orthography. This means that we need both a primary field and a Cyrillic field for a number of categories of information: lemma, ergative of nouns, plural of nouns, present tense of verbs, witnessed past of verbs. The alternative would be to put the orthography in the same field as the transcription, but this would mean that the dictionary cannot be sorted by the orthography.

Reference

Bell, John & Bird, Steven. 2000. A preliminary study of the structure of lexicon entries. Paper presented at the Workshop on Web-based Language Documentation and Description, 12-15 December 2000, Philadelphia (USA).