

Linking Basque Lexical Resources

Abstract

This poster presents preliminary considerations for a new project: A merged set of Basque (legacy) lexical resources or unified lexical database. At this preliminary stage, our main attention lies on the catalogue of data sources, on philological problems (e.g. regarding lemmatization), and on the design of the database. We propose a data model and a workflow for the inclusion of all kinds of Basque dictionaries and other resources such as the Basque WordNet and NLP lexicons. The data found in Basque dictionaries bare several problems, such as presence of dialectal and historical forms from before and after the creation of a Basque standard in 1968, and inconsistencies in lemmatization and the treatment of homography, homonymy, and polysemy. We present some solutions for these problems, as well as an XML schema for the lexical database and example datasets. The new merged resource may be used as Basque diachronic lexicographical corpus or serve as datasource for the creation of new lexicographical products. Using word sense identifiers found in Basque WordNet, the resource may also be linked to lexical resources of other languages.

Keywords: Historical Lexicography, Linking Lexical Resources, Basque Language

POSTER PROPOSAL – Language: English

1 Introduction

Two dates have a powerful meaning in Basque lexicography: On the one hand, the foundations for the creation of Standard Basque date from not earlier than 1968. On the other hand, from around 2005 Basque state-of-the-art electronic dictionaries and other resources have been published, and printed dictionaries are today fairly out of use and due to disappear. We propose to gather the data from Basque (legacy) lexicographical publications in an electronic database, with multiple goals in mind: (1) To preserve the lexicographical work done by numerous Basque scholars, (2) to make this work universally accessible, (3) to present all the data together, defining a common data model for this “dictionary of dictionaries”, which shall allow the linking of objects stemming from legacy dictionaries from before 1968, i.e. the macrostructure of which does not obey Standard Basque lemmatization, through (a) Standard Basque lemma signs, (b) Standard Basque lempos-entities (i.e., lemmata with a unique part of speech value assigned), and (c), WordNet concepts (i.e. word sense identifiers found in Basque WordNet (EUSWN, *cf.* Table 1 below).

2 The Catalogue

In Table 1, we list Basque dictionaries and lexical resources of other types.¹ The tables' four columns represent the following:

¹ The poster proposed here will contain sample pictures of paper dictionary content.

1. An abbreviation for each resource, typically consisting of a group of letters indicating title and/or authors;
2. The methods that according to preliminary assumptions apply for an inclusion of the data contained in the respective resource in the merged resource. The different methods are abbreviated as follows:
 - a) TP Transcription and Parsing (i.e., representation in an XML structured format)², for manuscripts and older printed material
 - b) OP Optical Character Recognition (OCR) and Parsing
 - c) P Parsing only, if the data are already accessible in an electronic format
 - d) X Transformations of the structural markup the data contains, for data already parsed or initially designed as electronic resource enriched with structural markup
3. An indication, whether the lemma list of the respective resource reflects nowadays Standard Basque (YES) or it does not or does so partially (NO).
4. A short version of the bibliographical reference.

At the moment, data from about a dozen of these sources is already accessible to us in an electronic format and has been imported in a preliminary version of the merged resource.

Dictionary	Method	Standard	Reference
ADO1980	OP	NO	Kintana et al. (1980): <i>Hiztegia 80: euskara-espainiera, espainiera-euskara, vasco-español, español-vasco</i>
ADO1984	P	NO	Kintana et al. (1984): <i>Hiztegia bi mila: euskara-espainiera espainiera-euskara, vasco-español español-vasco</i>
ADO1986	P	NO	Uribarren et al. (1986): <i>Euskararako Hiztegia</i>
ADO1993	P	NO	Uribarren et al. (1993): <i>Europa hiztegia: eskola berrirakoa</i>
ADO1996	P	YES	Uribarren et al. (1996): <i>3000 Hiztegia</i>
ADO2009	X	YES	Adorez Taldea (2009): <i>Adorez 5000 Hiztegia</i>
ADO2013	X	YES	Bostak Bat Kultur Elkartea (2013): 5000 Hiztegia .
AIZ1883	OP	NO	Aizquibel (1883). <i>Diccionario basco-español titulado euskeratik erderara bi-urtzeco itztegia</i>
AUWH1992	P	NO	Aulestia & White (1992): <i>Euskara-Ingelesa Ingelesa-Euskara Hiztegia</i> . Donostia: Elkar
AZK1905	OP	NO	Azkue (1984 [1905-06]). <i>Diccionario vasco-español-francés</i> . Bilbo: Euskaltzaindia [1905-06ko lehenengo argitalpenaren faksimilea]
BEME1916	OP	NO	Bera & Mendizabal (1916). <i>Diccionario castellano-euzkera / Bera'tar Eroman Mirena aba, buruñurduna. Euzkel-erdel iztegia / López Mendizabal'dar Ixaka</i>
BERA1909	OP	NO	Bera (1909). <i>Euzkel-iztegitxua / Bera'tar Erroman Aba</i>
CHA2002	X	NO	Charpentier (Charpentier 2002 [1823]): <i>Wörterverzeichnis, „in der zu Saint Etienne (Donostiä) Hauptdorf des Thales von Baigorri bei Saint-Jean-Pied-De-Port üblichen Mundart“</i>

² The authors have performed a series of tests for parsing the content of OEH dictionary using GROBID-dictionaries, a tool that relies on Conditional Random Fields, a Machine Learning approach, for parsing macrostructure and microstructural content, with promising results (cf. Blog Post [ANONYMIZED], 2017).

Dictionary	Method	Standard	Reference
EAH2007	P	YES	Martínez Rubio (2007): Euskara-Alemana Hiztegia.
EDBL	X	YES	Aldezabal et al. (2001) <i>Euskarazko Datu-Base Lexikala</i>
EEF2012	P	YES	Elhuyar Hiztegiak: Euskara-Frantsesa, 2. arg. (Pikabea et al. 2007) http://hiztegiak.elhuyar.org
EEG1996	P	YES	Elhuyar Hiztegiak: Euskara-Gaztelaina, 1. arg. (Azkarate et al. 1996)
EEG2000	P	YES	Elhuyar Hiztegiak: Euskara-Gaztelaina, 2. arg. (Azkarate et al. 2000)
EEG2006	T	YES	Elhuyar Hiztegiak: Euskara-Gaztelania, 3. arg. (Azkarate et al. 2006)
EEG2013	T	YES	Elhuyar Hiztegiak: Euskara-Gaztelaina, 4. arg. (Elhuyar Hizkuntza Zerbitzuak 2013) http://hiztegiak.elhuyar.org
EEN2007	T	YES	Elhuyar Hiztegiak: Euskara-Ingelesa, 2. arg. (Elhuyar Hizkuntza Zerbitzuak 2007) http://hiztegiak.elhuyar.org
EER1997	P	YES	Elhuyar Hiztegiak: Euskara-Errusiera, 1. arg. (Serrano et al. 1997)
EHA2012	X	YES	Euskaltzaindia (2012): <i>Euskaltzaindiaren hiztegia: Lesartk eta adibideak</i>
EM1987	OP	NO	Etxebarria, J.M. & Mujika, J.A. (1987). <i>Euskararen oinarriko hiztegia: maiztasun eta prestasun azterketa.</i>
ETX2001	P	YES	Etxeberria (2001): <i>Ikaslearen Hiztegia</i>
EUNL1996	OP	YES	Jansen, W.H. (1996): Baskisch-Nederlands, Nederlands-Baskisch = Euskara-nederlandera, nederlandera-euskara
EUSWN	X	YES	<i>Euskal WordNet</i> (Pociello 2007; Pociello et al. 2011)
HAR1741	TP	NO	in: Harriet (1741): <i>Gramatica escuaraz eta francesez, composatua francez hitzcunça ikhasi nahi dutenen faboretan.</i>
HB2000	P	YES	Euskaltzaindia (2000): <i>Hiztegi Batua</i>
HB2010	X	YES	Euskaltzaindia (2010): <i>Hiztegi Batua</i>
HUM1817	TP	NO	Humboldt (1817): <i>Berichtigungen und Zusätze...</i>
IKAS1976	OP	NO	Montiano & Urquijo (1976): <i>Diccionario "IKAS" Euskera-Castellano Castellano-Euskara. Dialectos vizcaíno y guipuzcoano</i>
KÜH1999	OP	NO	Kühnel (1999): <i>Wörterbuch des Baskischen</i>
LAR1745	TP	NO	Larramendi (1745): <i>Diccionario Trilingüe</i>
LECL1826	OP	NO	in: Lécuse (1826): <i>Grammaire Basque</i>
LHA1926	OP	NO	Lhande (1926). <i>Dictionnaire basque-français et français-basque: dialectes Labourdin, Bas-Navarraïns et Souletin</i>
LÖP1968	OP	NO	Löpelmann (1968): <i>Ethymologisches Wörterbuch der baskischen Sprache</i>
LUR1996	OP	YES	Arratibel, Atela & Navarro (1996): <i>LUR Eskolarako Hiztegi Entziklopedikoa</i>
MAH1840	TP	NO	Mahn (1840) [ms.]: Baskisches Wörterbuch I & II
MEND1962	OP	NO	López-Mendizabal (1962). <i>Diccionario vasco-español</i> , 4. arg.
MOD1977	OP	NO	Kintana (1977): <i>Euskal Hiztegi Modernoa</i>
MOD2000	P	YES	Barandiaran & Etxeberria (2000): <i>Euskal Hiztegi Modernoa</i> , 2. arg. berritua
MORR1998	X	YES	Morris (1998): Morris Student Plus. Euskara-Ingelesa, English-Basque
MUG1981	P	NO	Múgica Berrondo (1981): <i>Diccionario vasco-castellano</i>
OEH	X	NO	Euskaltzaindia (Mitxelena & Sarasola 1988): <i>Orotariko Euskal Hiztegia</i>
POUV1665	TP	NO	Pouvreau (1665) [ms.]: <i>Euskara-frantses(-latin-gaztelaniazko) hiztegia.</i>
RET1976	OP	NO	Sota, Lafitte & Akesolo (1976): <i>Diccionario Retana de autoridades de la lengua vasca: con cientos de miles de nuevas voces y acepciones, antiguas y modernas</i>
SEH1996	P	YES	Sarasola (1996): <i>Euskal Hiztegia</i>

Dictionary	Method	Standard	Reference
SEH2007	P	YES	Sarasola (2007): <i>Euskal Hiztegia</i>
SEIH1999	P	YES	Sarasola (1999): <i>Euskara Ikaslearen Hiztegia</i>
SHLH1994	X	YES	Sarasola (1994): <i>Hautu-lanerako Euskal Hiztegia</i> . http://www.euskara.euskadi.net/r59-sarasola/eu/sarasola/sarasola.apl
SMH1982	X	NO	Sarasola (1982). <i>Gaurko euskara idatziaren maiztasun-hiztegia: 1977ko corpus batean oinarritua</i>
UZEI1987	X	YES	Euskalterm (UZEI 1987): http://www.euskara.euskadi.net/r59-euskalte/eu/q91Eu-sTermWar/kontsultaJSP/q91aAction.do
UZEI2004	X	YES	UZEI (2004): <i>Maiztasun Hiztegia</i>
UZEI2006	X	YES	UZEI (2006): <i>Atzekoz aurrera. Hitz-bukaeren hiztegia</i>
ZEH2005	X	YES	Sarasola (2005): <i>Zehazki: gaztelania-euskara hiztegia</i> . http://www.ehu.eus/ehg/cgi/zehazki

Table 1: Basque lexical resources and methods for inclusion in a merged resource

3 Data Modelling

As further development of the proposal presented in ([ANONYMIZED] et al., 2016], we propose to build a single XML document³ with a DTD following the TEI-P5 guidelines. On the first hierarchy level, a list of Standard Basque lemma signs is set as macrostructure for the merged resource (TEI-P5 element <super-entry>). The entry-headwords from the different lexical resources to become part of the merged resource will be inserted as child elements hereof (TEI-P5 element <entry>). Below each of these, the respective microstructural content is to be inserted, with a microstructural markup as rich as possible, from case to case. Whenever possible, i.e. if the data of a particular source is accessible in or can be transformed to proper XML, we will establish a syntactical entity which carries the part-of-speech (POS) information and, at the same time, serves for disambiguating homographs and homonyms (i.e. homographs with the same POS), so that the other information contained in the dictionary articles (semantic information etc.) is stored inside an element that is child to that syntactical entity. This modelling is particularly interesting for Basque, since 5% of lemma signs have more than one part of speech assigned.

We are sure the Galway workshop will be a very good opportunity to discuss this approach of data modelling, and the technologies to be employed for retro-digitizing Basque legacy dictionaries, and to learn more about technologies available through or to be developed by the ELEXIS project.

³ The poster proposed here will contain figures representing the structure of that XML document.