



Computational Treatment of Basque Multiword Expressions

(WG1, WG2, WG3)

Ruben Urizar

IXA research group – University of the Basque Country (UPV / EHU)



Sources

EUSKERA
ESTADÍSTICA DE LA LINGÜÍSTICA
EUSKALTZAINdia
Hiztegi Batua

Hiztegi Batua dictionary (2010)

- prescriptive dictionary
- all subentries selected

Statistical Corpus of 20th Century Basque

- 4.7 million words
- balanced corpus
- manually tagged, including MWEs
- we selected all MWEs occurring 10 times or more

POS	
Verbs	837
Nouns	695
Adverbs	343
Quantifiers	113
Conjunctions	93
Adjectives	53
Interjections	33
Pronouns	20
Others	20
TOTAL	2,207

Estatistikak
Idazti galdetza, aukeratu
bikotako mota eta sareak
Bilatu, bultzatua

www.EuskaraCorpusa.net

Alaria XX. mendeko euskararen korpus estatistikoak

Testu-Aitzkariaren
Tarteano: Ordenean
Lurraldeko: Aitzkariaren

Aitzkariaren

1. 1959-1960 Euskarazko Saitz-aitzakuk C. Jensen 2005
2. Lan ospaini postzen du bizi (ezpido leinu...) 13 - Ezpido leinu buruko minak emanen ditzak.

3. 1959-1960 Salaketa gabeak Saitz-aitzakuk Ikermet 1975_0005
- On arteboroko lan bat egitea data

4. 1959-1960 Euskarazko Saitz-aitzakuk Suza 1959_0115
SUSTEN

5. 1959-1960 Euskarazko Saitz-aitzakuk J. Kortazar 2007
Ikermet Saitz-aitzakuk Praktikoa 1959-1960

6. 1959-1960 Euskarazko Saitz-aitzakuk X. Koltzschka 2014
Guru arteen oraindik kontzentru da zentzu hizkera edo aditz itzela kamarrek. Mbiluen harribit lan egin zuen hitz itzala edo "orriengrabi word" orriengrabi direkleku salatzen.

7. 1960-1961 Gipuzkoako Saitz-aitzakuk U. 0002
Argia (bulatenaren hiru buruan erabili diren aditz zamerak ez baute berarekin gatzera utzi, ez hizkera arrekekin) argia doanenak legeztatu.

Representation Lexical Database for Basque (EDBL)



- The purpose of the description is to formally encode all the possible surface realizations of each MWE.

- We worked out a single representation covering all types of MWEs ranging from fixed expressions to those of highest morphosyntactic flexibility.

The description of MWEs within EDBL includes, at least, three aspects:

- their **composition**, i.e. which the components of the MWE are, whether each of them can be **inflected** or not, and which one-word lexical unit conveys the **morphosyntactic information** to the whole MWE
- their **surface realization**, that is, the **order** in which the components may occur in the text, the mandatory or optional **contiguosity** of components, and the **inflectional restrictions** applicable to each one of the components.
Different realization patterns may be defined for each MWE
- their possible **ambiguity**, i.e. whether the sequence of words matching a given surface realization pattern must be unambiguously analyzed as an instance of the MWE or, on the contrary, may be analyzed as separate words in some contexts.

In total, we used 177 different realization patterns and 145 inflection restrictions.

Example (EDBL)

Composition of the MWE

COMPOSITION

- Lemma (sarrera): **aditzera eman** ‘to announce’ lit. ‘to give to understand’
- Part of Speech (kategoriak): Verb (ADI)
- Components (osagaia): VB **aditu** (understand) + VB **eman** (give)
- Component conveying morpho-syntactic information to the whole MWE: **eman**

SURFACE REALIZATION

For MWE **aditzera eman** ‘to announce’ there are 4 realization patterns corresponding to 4 possible orders of components:

- Order of components (orden- jarratzauna):
 - contiguous: 12 and 21
 - split: 1+2 and 2+1
- Inflectional restrictions (flexio-murriztapena):
 - first component (**aditzera**, ‘to understand’) is fixed [-]
 - second component (**eman**, ‘to give’) may take any inflection [+]
- Unambiguity (ziurra):
 - orders 12 and 21 are unambiguous
 - 1+2 and 2+1 are ambiguous.

MWE processor HABIL

- It deals with both **contiguous** and **split MWEs**
- It takes into account **all the possible orders** of the components
- It checks that **inflectional restrictions** are complied with
- It generates morphosyntactic **interpretations** for the MWE

Example (CG grammar rule)

RULE

```
ADD (%MWE)
TARGET LAN-EGIN
CONDITIONS (1P):
IF (0 EGIN)
  (1 EDUN/EZAN) (NOT 1 NOR-HAIKE)
  (2 LAN AND LAN-EGIN)
  (NOT 3 IZENONDO OR POSDET) -->
  0 position: 2nd component is a verb in list EGIN
  1 position: TR AUX (edun or ezan), OBJ NOT 3P
  2 position: 1st component is a verb in list LAN and marked as MWE in list LAN-EGIN
  3 position: NOT post.ADJ or post.DET
```

LISTS

```
LIST LAN-EGIN = "lan_egin" "hitz_egin" "ihes_egin" "parte_hartu" ;
LIST LAN = "lan" "hitz" "ihes" "parte" ;
LIST EGIN = "egin" "hartu" ;
```

```
LIST EDUN/EZAN = "*edun" "*ezan" ;
LIST NOR-HAIKE = "NOR_HAIKE" ;
LIST IZENONDO = (ADJ IZAUR-) ;
LIST POSDET = (DET ERKARR) (DET BAN)
  "berbera" "bat" "batzuk" "bi" "anitz" "aski" "asko" "dena"
  "franko" "guzti" "gutxi" "gehiago" "gehiiegia" ... ;
```

CG disambiguation grammar

- PREMISE:** Many MWEs may undergo more restricted variations than literal uses.
- We chose the **20 most frequent MWEs** described in the lexical database having at least one **ambiguous** realization pattern.
 - For the development of the grammar, we built a **sub-corpus of 21,125 sentences** from the **Statistical Corpus of 20th Century Basque**
 - The sub-corpus which contained occurrences of word combinations corresponding to both MWE and literal interpretations.
 - The grammar we have built consists of **111 rules**, many of which can be reused to disambiguate MWEs of the same type.
 - The grammar has proved to have 98.90% coverage and 98.88% accuracy.

CONCLUSION: For many Basque MWEs, morphosyntactic context can be enough to resolve ambiguity.

Raw text

Erlijio alorrean Zaitegik berak egin zuen lan batez ere.

religion area.INE Zaitegi.ERG himself.ERG do-PART AUX.3S.PST work above all

It was specially Zaitegi himself who worked on the religion area.

Morphosyntactical analysis: EUSTAGGER

```
<<Erlijio>><HAS MAI>
"erlijio" IZE ARR ZERO HAS MAI @KM>
"erlijio" IZE ARR ABS MG HAS MAI @OBJ @PRED @SUBJ
<<alorrean>
"alor" IZE ARR INE NUMS MUGM @ADLG
<<Zaitegik>><HAS MAI>
"Zaitegi" IZE LIB ERG NUMS MUGM HAS MAI @SUBJ
<<berak>
"bera" ADJ ARR ERG MG AORG @SUBJ
"bera" ADJ ARR ABS NUMS MUGM AORG @OBJ @PRED @SUBJ
"bera" DET ERKIND NMGP ABS NUMP MUGM @OBJ @PRED @SUBJ
"bera" DET ERKIND NMGS ERG NUMS MUGM AORG @SUBJ
<<egin>
"egin" ADI SIN ADDIN NOTDEK @JADNAG
"egin" ADI SIN PART BURU NOTDEK @JADNAG
"egin" IZE ARR ABS MG @OBJ @PRED @SUBJ
[...]
<<zuen>
"edun" ADL B1 NOR NORK NR HURA NK HARK @JADLAG
"ukan" ADT PNT B1 NOR NORK NR HURA NK HARK @JADNAG
"zuek" IOR PERARR ZUEK GEN NUMP MUGM ZERO @IZLG @IZLG>
[...]
<<lan>
"landu" ADI SIN ADDIN NOTDEK @JADNAG
"lan" IZE ARR ABS MG @OBJ @PRED @SUBJ
[...]
<<batez>
"bate" IZE ARR INS MUG @ADLG
"bat" DET DZG INS MUG @ADLG
"batez" IZK ZKI INS MUG @ADLG
<<ere>
"ere" LOT LOK EMEN @LOK
"<$,><PUNT PUNT>" PUNT PUNT
```

MWEs detected: HABIL

```
<<Erlijio>><HAS MAI>
"erlijio" IZE ARR ZERO HAS MAI @KM>
"erlijio" IZE ARR ABS MG HAS MAI @OBJ @PRED @SUBJ
<<alorrean>
"alor" IZE ARR INE NUMS MUGM @ADLG
<<Zaitegi-k>><HAS MAI>
"Zaitegi" IZE LIB ERG NUMS MUGM HAS MAI @SUBJ
<<berak>
"bera" ADJ ARR ERG MG AORG @SUBJ
"bera" ADJ ARR ABS NUMS MUGM AORG @OBJ @PRED @SUBJ
"bera" DET ERKIND NMGP ABS NUMP MUGM @OBJ @PRED @SUBJ
"bera" DET ERKIND NMGS ERG NUMS MUGM AORG @SUBJ
<<egin>><1-2>
"egin" ADI SIN ADDIN NOTDEK @JADNAG
"egin" ADI SIN PART BURU NOTDEK @JADNAG
"egin" IZE ARR ABS MG @OBJ @PRED @SUBJ
"lan_egin" ADI ADK ADDIN NOTDEK mw1 @JADNAG
"lan_egin" ADI ADK PART BURU NOTDEK mw1 @JADNAG
[...]
<<zuen>
"edun" ADL B1 NOR NORK NR HURA NK HARK @JADLAG
"ukan" ADT PNT B1 NOR NORK NR HURA NK HARK @JADNAG
"lan_egin" ADI ADK ADDIN NOTDEK mw1 @JADNAG
"lan_egin" ADI ADK PART BURU NOTDEK mw1 @JADNAG
[...]
<<batez_ere>
"batez_ere" LOT LOK EMEN mw2 @LOK
"<$,><PUNT PUNT>" PUNT PUNT
```

MWEs disambiguated: CG grammar

```
<<Erlijio>><HAS MAI>
"erlijio" IZE ARR ZERO HAS MAI @KM>
"erlijio" IZE ARR ABS MG HAS MAI @OBJ @PRED @SUBJ
<<alorrean>
"alor" IZE ARR INE NUMS MUGM @ADLG
<<Zaitegi-k>><HAS MAI>
"Zaitegi" IZE LIB ERG NUMS MUGM HAS MAI @SUBJ
<<berak>
"bera" ADJ ARR ERG MG AORG @SUBJ
"bera" ADJ ARR ABS NUMS MUGM AORG @OBJ @PRED @SUBJ
"bera" DET ERKIND NMGP ABS NUMP MUGM @OBJ @PRED @SUBJ
"bera" DET ERKIND NMGS ERG NUMS MUGM AORG @SUBJ
<<egin>><1-2>
"egin" ADI SIN ADDIN NOTDEK @JADNAG
"lan_egin" ADI ADK ADDIN NOTDEK mw1 @JADNAG
"lan_egin" ADI ADK PART BURU NOTDEK mw1 @JADNAG
[...]
<<zuen>
"edun" ADL B1 NOR NORK NR HURA NK HARK @JADLAG
"ukan" ADT PNT B1 NOR NORK NR HURA NK HARK @JADNAG
"lan_egin" ADI ADK ADDIN mw1 @JADNAG
"lan_egin" ADI ADK PART BURU NOTDEK mw1 @JADNAG
[...]
<<batez_ere>
"batez_ere" LOT LOK EMEN mw2 @LOK
"<$,><PUNT PUNT>" PUNT PUNT
```

Ambiguity resolved