

Formalizing the Morphology of Rromani Adjectives

**Masako Watabe, University of Franche-Comté, CRIT, UFR SLHS,
30-32 rue Mégevand, F-25000 Besançon, France, masako.watabe@univ-fcomte.fr**
**Max Silberztein, University of Franche-Comté, CRIT, UFR SLHS,
30-32 rue Mégevand, F-25000 Besançon, France, max.silberztein@univ-fcomte.fr**

Abstract

This paper presents a set of linguistic resources that formalizes the morphological behavior of simple Rromani adjectives. We describe the formalization of the adjectives' morphology and the implementation with the NooJ linguistic platform of an electronic dictionary associated with a formal morpho-syntactic grammar. We can then apply this set of resources to a corpus to evaluate the resources and automatically annotate adjectival forms in Rromani texts. The final set of resources can then be used to identify each Rromani dialectal variant and can be used as a pedagogical tool to teach Rromani as a second language.

1 Introduction

1.1 Rromani language

Rromani is the language of the Rromani people; it is an Indo-Aryan language. The number of Rromani speakers is estimated at 5.5 million (Gurbetovski, M. et al. 2010). UNESCO's "Atlas of the World's Languages in Danger" classifies Rromani as a "definitely endangered¹" language (UNESCO. 2010). There are four Rromani dialects, formed by two isoglosses combining with each other (Courthiade, M. 2016):

- The first isoglossal criterion concerns the opposition between "o" and "e," e.g., *phirdom* vs. *phirdem* [I walked], *o Rroma* vs. *e Rroma* [the Rroms].
- The second isoglossal criterion concerns the phonetic mutation of two consonants: the alveolar affricates "ʃ^h" and "dʒ" transform

into alveolar-palatal fricatives [ç] and [ʒ], e.g., "ʃ^havo" vs. "çavo" [Rromani boy, son], "dʒukel" vs. "ʒukel" [dog].

These four dialects are not areal: Rromani speakers living in nearby regions do not necessarily speak the same dialects, and the same dialect is used in distant countries.

The Rromani alphabet was standardized at the International Rromani Union Congress in 1990, see Figure 1.

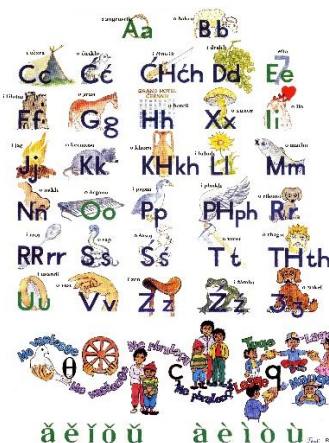
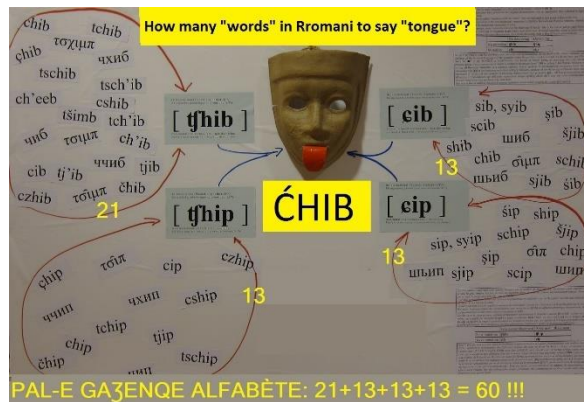


Figure 1: The Rromani standardized alphabet

If all Rromani speakers transcribe, for example, the word *çhib* [language] using their local alphabets, there can be up to 60 different spellings. The written word *çhib* is an underlying form including four possible pronunciations: [ʃ^hb], [ʃ^hp], [çib], and [çip], see Figure 2. The standardized alphabet enables speakers of different dialects to understand each other in writing, giving them comfort in pronunciation.

No other standardization exists: neither lexical, nor grammatical, nor phonetic.

¹ The UNESCO list has six categories of danger: Stable yet threatened, vulnerable, definitely endangered, severely endangered, critically endangered, extinct.



⁴ French is the default language of the Facebook account of one of the authors, *i.e.*, Watabe, M.

MS.) and a one-page poem (Đurić, R. 2006), contained only 747 lexical entries associated with a well-developed morphological grammar that includes 179 inflectional paradigms and 11 derivational paradigms⁵ for nouns, verbs, adjectives, and grammatical words. A feature of these resources is that they take into account Rromani four dialects, as well as a few vernaculars.

An editorial dictionary (Courthiade, M. et al. 2009) including the four dialects of Rromani explains Rromani morphology in the grammar section. It is our principal lexical and grammatical resource.

In this paper, we are addressing the problem of describing adjectives and their inflection, which causes massive ambiguities.

3 Rromani adjectives

The inflectional morphology of Rromani adjectives is governed by two genders (masculine and feminine), two numbers (singular and plural), and two cases (direct and oblique⁶). Adjectival forms are according to noun genders, numbers, and cases. The basic form of adjectives is the masculine singular direct. Combining these three properties produces eight possibilities; in practice, however, most adjectives have no more than three forms (Courthiade, M. et al. 2009. Sarău, G. 2009). Consequently, there are many inflectional homonyms.

Most Rromani words are oxytonic; *i.e.*, the tonic stress is on the last syllable. One uses a grave accent to mark the stress when it is not on the last syllable. For example, *bakri* [ewe] and *thulo* [fat, thick, dense] are oxytonic, whereas *profesòri* [professor] and *sociàlo* [social] are non-oxytonic. The opposition oxytonic *vs.* non-oxytonic plays a role in inflectional morphology.

3.1 Oxytonic adjectives

Oxytonic adjectives are classified into four classes: large adjectives, narrow adjectives, plural adjectives, and invariable adjectives.

Large adjectives⁷: Large adjectives are vocalic and have three suffixes: “-o” in the masculine singular direct, “-i” in the feminine singular direct, and “-e” in the plural direct for both genders, as well as oblique for both genders and numbers, see an example of the adjective *thulo* [fat, thick, dense] in Table 1:

Form	Gender	Number	Case
<i>thulo</i>	masculine	singular	direct
<i>thuli</i>	feminine	singular	direct
<i>thule</i>	masculine	plural	direct
<i>thule</i>	feminine	plural	direct
<i>thule</i>	masculine	singular	oblique
<i>thule</i>	feminine	singular	oblique
<i>thule</i>	masculine	plural	oblique
<i>thule</i>	feminine	plural	oblique

Table 1: Inflected forms and properties of the adjective *thulo* [fat, thick, dense]

The form *thule* is therefore 6-time ambiguous. This high level of ambiguity is general in Rromani; as a matter of fact, we do not know any Rromani adjective that would inflect to eight different forms, each for each combination of properties.

Narrow adjectives⁸: Narrow adjectives are consonant, and the direct forms of both genders are identical in each number: “-Ø” in the direct singular and “-a” in the direct plural. The suffix of the oblique is “-e” for both genders and numbers, as in “large” adjectives, see an example of the adjective *godăver* [intelligent] in Table 2:

Form	Gender	Number	Case
<i>godăver</i>	masculine	singular	direct
<i>godăver</i>	feminine	singular	direct
<i>godăvera</i>	masculine	plural	direct
<i>godăvera</i>	feminine	plural	direct
<i>godăvere</i>	masculine	singular	oblique
<i>godăvere</i>	feminine	singular	oblique
<i>godăvere</i>	masculine	plural	oblique
<i>godăvere</i>	feminine	plural	oblique

Table 2: Inflected forms and properties of the adjective *godăver* [intelligent]

⁵ Only the diminutive and abstract nouns with the suffix “-pen” are described in the current Rromani module.

⁶ In Rromani, the direct case of human and most animal nouns is used as a subject, while the oblique case is used as an object complement. The direct case of inanimate object nouns is used as a subject and an object complement.

⁷ The Rromani adjective *buxlo* [large] is the origin of this designation.

⁸ The Rromani adjective *tang* [narrow] is the origin of this designation.

Plural adjectives: Plural adjectives are used specifically with plural nouns. In the direct case, they have the suffix “-Ø” for both genders, and in the oblique “-e” for both genders, see an example of the adjective *but* [many, numerous] in Table 3:

Form	Gender	Number	Case
<i>but</i>	both	plural	direct
<i>bute</i>	both	plural	oblique

Table 3: Inflected forms and properties of the adjective *but* [many, numerous]

Invariable adjectives: So-called “international” adjectives have a tendency to be invariable. International oxytonic adjectives are completely invariable, see an example of the adjective *bordo* [Bordeaux-colored] in Table 4:

Form	Gender	Number	Case
<i>bordo</i>	both	both	both

Table 4: Invariable form and properties of the adjective *bordo* [Bordeaux-colored]

3.2 Non-oxytonic adjectives

Borrowed adjectives and suffixed adjectives:

Inflectional paradigms of borrowed non-oxytonic adjectives and suffixed non-oxytonic adjectives are identical. Their suffix is “-o” in the direct singular for both genders, “-a” in the direct plural for both genders, and “-one” in the oblique for both genders and numbers. Oblique forms are oxytonic, therefore the stress is not marked, see an example of the adjective *vešitko* [of the woods, wild] in Table 5:

Form	Gender	Number	Case
<i>vešitko</i>	masculine	singular	direct
<i>vešitko</i>	feminine	singular	direct
<i>vešitka</i>	masculine	plural	direct
<i>vešitka</i>	feminine	plural	direct
<i>vešitkone</i>	masculine	singular	oblique
<i>vešitkone</i>	feminine	singular	oblique
<i>vešitkone</i>	masculine	plural	oblique
<i>vešitkone</i>	feminine	plural	oblique

Table 5: Inflected forms and properties of the adjective *vešitko* [of the woods, wild]

International adjectives: Compared to international oxytonic adjectives (e.g., *bordo* [Bordeaux-colored]), international non-oxytonic adjectives are not completely invariable. International non-oxytonic adjectives have two suffixes: “-o” in the direct for both genders and

numbers and “-one” in the oblique for both genders and numbers. Oblique forms are oxytonic, therefore the stress is not marked, see an example of the adjective *sociàlo* [social] in Table 6:

Form	Gender	Number	Case
<i>sociàlo</i>	masculine	singular	direct
<i>sociàlo</i>	feminine	singular	direct
<i>sociàlo</i>	masculine	plural	direct
<i>sociàlo</i>	feminine	plural	direct
<i>socialone</i>	masculine	singular	oblique
<i>socialone</i>	feminine	singular	oblique
<i>socialone</i>	masculine	plural	oblique
<i>socialone</i>	feminine	plural	oblique

Table 6: Inflected forms and properties of the adjective *sociàlo* [social]

3.3 Conclusion

We have defined six classes of Rromani adjectives, according to their morphological properties:

- Oxytonic vocalic adjectives ending in “-o”: e.g., *thulo* [fat, thick, dense], *buxlo* [large],
- Oxytonic consonant adjectives: e.g., *godäver* [intelligent], *tang* [narrow],
- Oxytonic consonant adjectives used only in the plural: e.g., *but* [many, numerous],
- Oxytonic vocalic international adjectives totally invariable: e.g., *bordo* [Bordeaux-colored], *pane* [breaded],
- Non-oxytonic borrowed adjectives and non-oxytonic suffixed adjectives: e.g., *zèleno* [green], *vešitko* [of the woods, wild],
- Non-oxytonic vocalic international adjectives ending in “-o”: e.g., *sociàlo* [social], *interesànto* [interesting].

4 Formalization of the Rromani vocabulary

4.1 The NooJ linguistic platform

NooJ is a linguistic development environment linguists use to describe natural languages, by constructing linguistic resources in the form of electronic dictionaries and formal grammars from the Chomsky-Schützenberger hierarchy: regular, context-free, context-sensitive, and unrestricted

grammars. NooJ can formalize twelve levels of linguistic phenomena, from the typographical to the semantic level (Silberstein, M. 2016).

To formalize the Rromani adjectives vocabulary, one needs to construct the following linguistic resources:

- a dictionary containing the affixes, simple words, compound words, and discontinuous expressions that make up the Rromani vocabulary of adjectives
- a grammar containing the description of adjectives inflectional paradigms

One could describe Rromani's four dialectal variants in the dictionary and morphological grammar. The following sections present these levels of description.

4.2 Electronic dictionary

For example, the adjective *thulo* is represented by the following lexical entry in NooJ formalized (aka electronic) dictionary:

thulo,ADJ+EN="fat, thick, dense"+FLX=BUXLO
+DRV=ĆACĪPEN:SASTIPEN

In this extract, each lexical entry is composed of a lemma, its category "ADJ" (adjective), its English translation "+EN," and the name of its inflectional paradigm "+FLX".

The lexical entry *thulo* is associated with derivational paradigm ĆACĪPEN and its derivative's inflectional paradigm SASTIPEN. ĆACĪPEN describes the derivation of abstract nouns with the suffix "-pen," which applies to words of various categories, such as adjectives, nouns, and verbs.

The four main dialects are encoded using the following double codes:

- O-bi dialect: rro+rrbi
- O-mu dialect: rro+rrmu
- E-bi dialect: rre+rrbi
- E-mu dialect: rre+rrmu

In addition, we have added a third code to label specific language variants. For example, the

northern speech used in Russia and Poland is defined by the extra label: "+rrn." This is the case of the entry *vešitko* [of the woods, wild].

vešitko,ADJ+rro+rrbi+rrn+EN="of the woods, wild"+FLX=VEŠĪTKO+SYN="vešutno"

The entry above shows that the adjective *vešitko* belongs to the dialects O-bi (+rro+rrbi), and is used specifically in Russia and Poland (+rrn), its English translation is "of the woods, wild" (+EN="of the woods, wild"), it is inflected according to the paradigm named VEŠĪTKO (+FLX=VEŠĪTKO⁹), and it has the synonym "vešutno" (+SYN=vešutno) used in most Rromani dialects except the vernacular in Russia and Poland.

4.3 NooJ morphological grammar

In NooJ, inflectional paradigms are represented by regular or context-free grammars built over suffix/property factors: suffixes are added to the lexical entry to construct forms, which are associated with the corresponding properties (Silberstein, M. 2003-). For example, the following is the grammar rule that describes the inflectional paradigm RROM:

RROM = <E>/sg+dr | a/pl+dr | es/sg+ob | en/pl+ob;

This rule states that if one adds the empty string (<E>) to the lexical entry, one produces a singular (+sg) direct (+dr) form; if one adds an "a" to the lexical entry, one produces a plural (+pl) direct (+dr) form; if one adds "es" to the lexical entry, one produces a singular (+sg) oblique (+ob) form; if one adds "en" to the lexical entry, one produces a plural (+pl) oblique (+ob) form.

Suffixes may contain stack operators. For instance, operator (for "Backspace") is used to delete the current letter. In the following paradigm:

BUXLO = <E>/m+sg+dr | i/f+sg+dr |
e/m+pl+dr | e/f+pl+dr | e/m+sg+ob |
e/f+sg+ob | e/m+pl+ob | e/f+pl+ob |
:CMP/comparative ;

The second term states that if one deletes the last letter of a lexical entry and then adds an "i" (suffix

⁹ The lexical entry (i.e., *vešitko*) is in lower case, whereas the paradigm name (i.e., *VEŠĪTKO*) is in upper case. It

prevents confusion between two distinct values represented by identical writing.

i), one produces the feminine, singular, direct form (f+sg+dr) of the lexical entry.

For example, when the paradigm BUXLO is applied to the lexical entry *thulo* [fat, thick, dense], there will be no change (<E>) to the direct masculine singular, one final letter will be deleted () and an “i” will be added to the direct feminine singular, and one final letter will be deleted and an “e” will be added to produce the direct plural forms in both genders, and to the oblique forms in both genders and numbers. This paradigm then represents the three forms of *thulo*:

- *thulo*: masculine singular direct
- *thuli*: feminine singular direct
- *thule*: masculine plural direct, feminine plural direct, masculine singular oblique, feminine singular oblique, masculine plural oblique, or feminine plural oblique

It means that the wordform *thule* is associated with six potential linguistic analyses.

The last term of the BUXLO paradigm is used to produce the comparative forms of the lexical entries. :CMP refers to the name of the following rule (similarly to auxiliary symbols in generative context-free grammars):

CMP = eder/m+sg+dr | eder/f+sg+dr |
 edera/m+pl+dr | edera/f+pl+dr | edere/m+sg+ob |
 edere/f+sg+ob | edere/m+pl+ob | edere/f+pl+ob ;

The comparative suffix “-eder” is added in place of the final letter “o” in the *thulo* lexical entry (see “:CMP” in the paradigm BUXLO above). The comparative is declined like the narrow adjectives: without suffix in the direct singular of both genders, “-a” in the direct plural of both genders, and “-e” in the oblique of both genders and numbers. This rule produces three forms: *thuleder*, *thuledera*, and *thuledere* for eight linguistic analyses.

Beside a dozen generic operators such as that are available for any language, NooJ offers linguists the possibility of creating specific operators for each language. For instance, the Spanish operator <Á> is used to add an acute accent to the current vowel; the Hebrew operator <F> is used to de-finalize the last consonant of a word; the Tamazight operator <T> replaces letter “ḍ” with “ṭ”, etc. For the Rromani language, we have implemented two specific operators:

- <A> deletes a grave accent, regardless the position and returns to the initial position,
- <À> adds a grave accent to the current letter.

For example, operator <A> is used in the following paradigms:

VEŠÌTKO = <E>/m+sg+dr | <E>/f+sg+dr |
 a/m+pl+dr | a/f+pl+dr | <A>ne/m+sg+ob |
 <A>ne/f+sg+ob | <A>ne/m+pl+ob |
 <A>ne/f+pl+ob ;

In the VEŠÌTKO paradigm, the fifth term states that if one removes the grave accent of the lexical entry, and then adds the suffix “ne” (<A>ne), one produces the masculine singular oblique form (+m+sg+ob) of the lexical entry. The operator <A> is typically used in paradigms associated with non-oxytonic words. For example, if the paradigm VEŠÌTKO is applied to the lexical entry *zèleno* [green], its oblique form will be *zelenone*, i.e., without the grave accent.

By applying all inflectional NooJ paradigms to the dictionary, NooJ produces all the inflected forms for each lexical entry automatically. When applying these resources to a text, NooJ annotates all recognized word forms. For example, the wordform *vešitkone* will be annotated as an adjective (ADJ), its basic form is *vešitko*, its inflectional value is the oblique (ob), its dialect value is a Northern vernacular belonging to the O-bi dialect (rro+rrbi+rrn), its synonym in other dialects is *vešutno*. That helps users pedagogically recognize the relationship between the basic form, an inflected form, and a dialectal variant. However, there are four sets of annotations because oblique forms are identical for both genders and numbers, see Figure 6. We need syntactic grammar to resolve ambiguities.

It is often better for pedagogical applications, to use NooJ graphical grammars to describe some phenomena. For instance, the paradigm BUXLO shown above can equally be described with the following graph, see Figure 7.

vešitko	vešitko
0	
vešitko,ADJ+SuperDict=rr+Mutation=rrbi+Vernacular=rrn+EN=	
vešitko,ADJ+SuperDict=rr+Mutation=rrbi+Vernacular=rrn+EN=	
vešitko,ADJ+SuperDict=rr+Mutation=rrbi+Vernacular=rrn+EN=	
vešitko,ADJ+SuperDict=rr+Mutation=rrbi+Vernacular=rrn+EN=	
vešitko,ADJ+SuperDict=rr+Mutation=rrbi+Vernacular=rrn+EN=	
+EN="of the woods, wild"+SYN="vešutno"+Cs=ob+Gd=f+Nb=sg	
+EN="of the woods, wild"+SYN="vešutno"+Cs=ob+Gd=f+Nb=pl	
+EN="of the woods, wild"+SYN="vešutno"+Cs=ob+Gd=m+Nb=sg	
+EN="of the woods, wild"+SYN="vešutno"+Cs=ob+Gd=m+Nb=pl	

Figure 6: Inflected form *vešitkone* annotated automatically

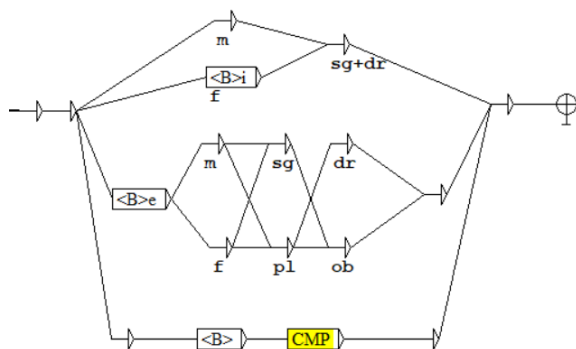


Figure 7: Inflectional grammar for adjective paradigm BUXLO.

The yellow node CMP refers to the embedded graph that represents the paradigm of the same name, see Figure 8.

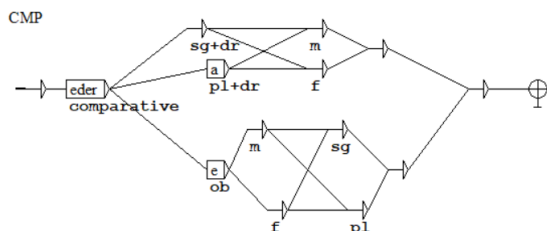


Figure 8: Inflectional grammar of the comparative suffix CMP.

The “DRV” property allows us to describe derivations. For example, the lexical entry *thulo* mentioned above:

thulo,ADJ+EN="fat, thick, dense"+FLX=BUXLO
+DRV=ĆACÍPEN:SASTIPEN

states that the word *thulo* can be derived according to paradigm $\acute{C}\acute{A}\acute{C}\acute{I}\acute{P}\acute{E}\acute{N}$, and the resulting derived form can be inflected according to inflectional paradigm $S\acute{A}S\acute{T}\acute{I}\acute{P}\acute{E}\acute{N}$. Following is the definition of the derivational paradigm $\acute{C}\acute{A}\acute{C}\acute{I}\acute{P}\acute{E}\acute{N}$:

$$\acute{C}\acute{A}\acute{C}IPEN = \langle B \rangle \langle A \rangle ipen/N+ina+m+abstract ;$$

From the lexical entry *thulo*, one deletes the final letter (), removes the accent (<A>¹⁰), and then adds suffix “ipen.” The resulting derivative *thulipen* [fatness, thickness, density], is a masculine inanimate noun (N+ina+m), that belongs to the semantic class “abstract.”

The abstract noun *thulipen* is a generic masculine singular direct form, and has seven dialectal variants: *thulipe* used in the two dialects O-bi and E-bi belonging to a dialectal subgroup “without mutation (rrbi),” *thulipo* used in the O-mu dialect (rro+rrmu), *thulimos* and *thulimo* used in the E-mu dialect (rre+rrmu), *thuliben*, *thulibe*, and *thulibo* used in the Carpathian vernacular belonging to the O-bi dialect (rro+rrbi+rrc), all of which mean “fatness, thickness, density.”

All these variants share the same inflectional morphology; *thulimàta* is the masculine plural direct form, *thulimas-* and *thulipnas-*¹¹ are the masculine singular oblique forms, and *thulimaten-* is the masculine plural oblique form. The following grammar rule describes the inflectional forms of the derivative as well as all its dialectal variants:

SASTIPEN = <E>/sg+dr | <B3>màta/pl+dr |
 <B3>(mas|pnas)<E>/sg+ob |
 <B3>maten<E>/pl+ob |
 /sg+dr+rrbi | <B2>o/sg+dr+rr+rrmu |
 <B3>(mos|mo)<E>/sg+dr+rre+rrmu |
 <B3>(ben|be|bo)<E>/sg+dr+rr+rrbi+rrc ;

NooJ automatically produces the derived forms for each lexical entry, and annotates all recognized derived forms in texts. For example,

¹⁰ The <A> operator concerns only non-oxytonic words. Thus, it produces no change for the lexical entry *rromano*.

¹¹ This is an archaic form.

thulpe

0

thulo,N+INITIALCATEGORY=ADJ+EN="dense, fat, thick"+Nb=

'+Nb=sg+C5=dr+Mutation=nbi+Spec=ina+Gd=m+abstract

4.4 Automatic Natural Language Processing

NooJ can use the same linguistic resources both to parse and generate texts. For example, one can apply a dictionary and its corresponding inflectional grammar to automatically produce all the forms associated with each lexical entry of the dictionary, see an extract in Figure 10.

```
thuledera,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thuledera,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACI
thuleder,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thule,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thule,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thule,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thule,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thule,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thule,thulo,ADJ+EN="dense, fat, thick"+FLX=BUXLO+DRV=CACIPI
thulipen,thulo,N+INITIALCATEGORY=ADJ+EN="dense, fat, thick"+FL
thulipe,thulo,N+INITIALCATEGORY=ADJ+EN="dense, fat, thick"+FL
thulipo,thulo,N+INITIALCATEGORY=ADJ+EN="dense, fat, thick"+FL
thuliblen,thulo,N+INITIALCATEGORY=ADJ+EN="dense, fat, thick"+FL

.DRV=CACIPIEN:SASTIPEN+comparative+tpl+dr+m
.DRV=CACIPIEN:SASTIPEN+comparative+tpl+dr+f
DRV=CACIPIEN:SASTIPEN+comparative+sg+dr+m
DRV=CACIPIEN:SASTIPEN+comparative+sg+dr+f
CACIPIEN:SASTIPEN+f+pl+dr
CACIPIEN:SASTIPEN+f+pl+ob
CACIPIEN:SASTIPEN+f+sg+ob
CACIPIEN:SASTIPEN+m+pl+dr
CACIPIEN:SASTIPEN+m+pl+ob
CACIPIEN:SASTIPEN+m+sg+ob
CACIPIEN:SASTIPEN+m+sg+dr
thick"+FLX=BUXLO+DRV=CACIPIEN:SASTIPEN+sg+dr+ina+m+abstract
thick"+FLX=BUXLO+DRV=CACIPIEN:SASTIPEN+sg+dr+rri+ina+m+abstract
thick"+FLX=BUXLO+DRV=CACIPIEN:SASTIPEN+sg+dr+rrr+ina+m+abstract
thick"+FLX=BUXLO+DRV=CACIPIEN:SASTIPEN+sg+dr+trr+rri+ina+m+abstract
```

NooJ uses the same resources to parse texts, lemmatize and annotate their wordforms, to apply queries in the form of regular, context-free, context-sensitive, or unrestricted grammars, perform statistical analyses, compute semantic analyses in Predicative or XML format,

Locate a pattern in Jeta2016

Pattern is:

- ☐ a string of characters:
- ☐ a PERL regular expression:
- ☒ a Nool regular expression:

<but>

☐ a Nool grammar:

Set

☐ Syntactic Analysis

Index

- ☐ Shortest matches
- ☒ Longest matches
- ☐ All matches

Limitation

- ☐ All occurrences
- ☒ Only: 100 occ.
- ☐ 1 occ. per match

☒ Reset Concordance

11 12 13 14

Figure 11: The query “<but>” and its resulting concordance.

The wordform *but* corresponds to three semantic values: the adjective “many, numerous,” e.g., *phirdòmsa bute Themenθe* [I traveled to many countries], the adverb “a lot, much,” e.g., *but dikhlom* [I have seen a lot], and another adverb “very,” e.g., *but súkar siklile on* [They learned very well].

Syntactic grammar is underdeveloped in the current NooJ module for Rromani, this is why wordforms remain ambiguous in general.

However, the query “<but,ADV>” will not retrieve the adjectival inflected forms *bute* because adverbs are invariable, see Figure 12.

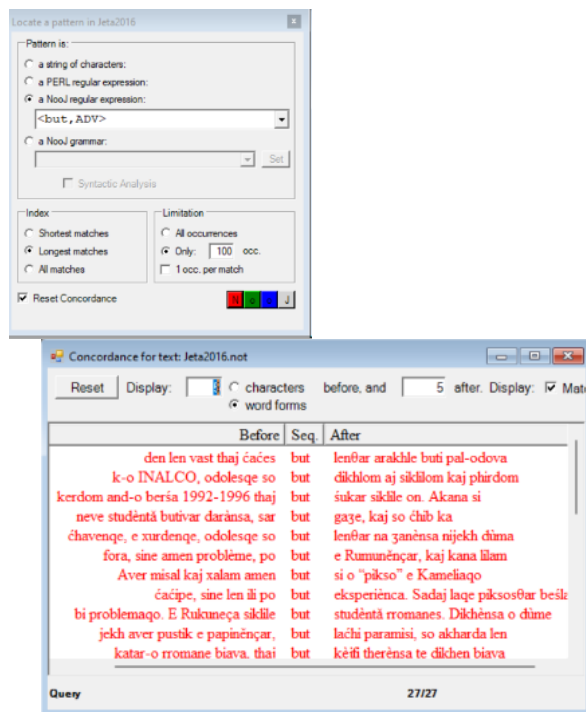


Figure 12: The query “<but,ADV>” and its resulting concordance.

If NooJ recognizes ambiguous forms, NooJ will show the annotation of all of them and not choose one randomly, as often happens in empirical applications. For example, as mentioned above, the adjectival inflected form *thule* is ambiguous because of six inflectional homonyms and the nominal inflected form *bakră* is ambiguous because of two inflectional homonyms, see Figure 13.

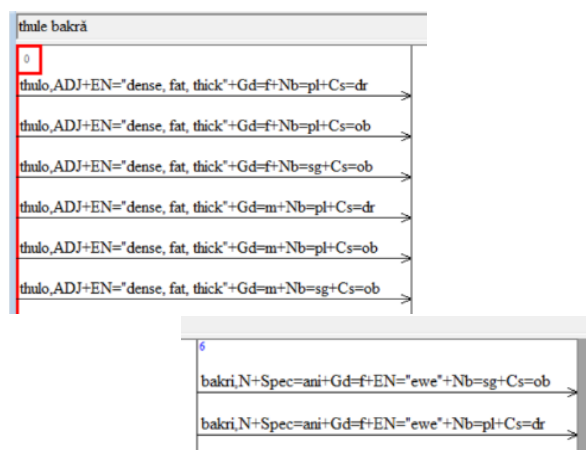


Figure 13: Inflected forms *thule* and *bakră* annotated automatically with ambiguity

4.5 Evaluation

There are 88 adjectives out of 747 lexical entries in the current NooJ dictionary for Rromani.

Applying these lexical entries and their corresponding inflectional grammars generated 1,278 inflected and derivational forms.

In our corpus, all wordforms that might correspond to potential adjectives have been recognized and annotated correctly, *i.e.*, we have reached a 100% recall, which is expected as we are specifically constructing our linguistic resources from the corpus. However, without any syntactic grammar, wordforms that might function as adjectives or as adverbs (*e.g.*, the wordform *but*) and be associated with different properties (*e.g.*, the wordforms *thule* and *bakră*) remain ambiguous until we can apply a syntactic grammar.

5 Conclusion, perspective

The current Rromani module recognizes all 170 adjectival forms from a small corpus that contains 708 wordforms. We are currently importing around 4,500 lexical entries from an editorial dictionary (Courthiade, M. et al. 2009) into a formalized NooJ format.

Removing ambiguities is our current challenge. We are constructing syntactic local grammars to disambiguate frequent adjectives.

The resulting linguistic resources will be downloadable from the NooJ website. The NooJ dictionary for Rromani will use the standard Rromani alphabet and include dialectal variants at the lexical and morphological levels. It will be available as a new digital and linguistic tool for all speakers of Rromani: native speakers and learners, regardless of their dialects.

We believe this polylectal dictionary is valuable from a dialectological point of view. Furthermore, as the declaration of the first Congress of the International Rromani Union in 1971 stated that “no dialect is better than another,” the dictionary will describe all dialects.

Unlike empirical methods, the NooJ platform produces analyses based on handcrafted linguistic resources, and thus offers linguists to describe and understand its properties. We believe that carefully and precisely handcrafting linguistic resources for Rromani is a worth scientific project, and will have many applications in Natural Language Processing, second-language teaching and corpus linguistics.

References

- Gheorghe Sarău. 2009. *Strukturë rromane çhibăqe*. Editura Universității din București, Bucharest.
- Jeta Duka. Deś berś vaś-i rromani çhib and-o INALCO. MS.
- Marcel Courthiade. Structure dialectale de la langue rromani. *Études tsiganes*, 22-2005, pages 14-26. Le Centre de documentation, Paris.
- Marcel Courthiade. 2016. The nominal flexion in Rromani. In Marcel Courthiade and Delia Grigore (eds.) *Professor Gheorghe Sarău: a life devoted to the Rromani language*. pages 157-211. Editura Universității din București, Bucharest.
- Marcel Courthiade et al. 2009. *Morri anghuni rromane çhibăqi evroputni lavustik*. Romano Kher, Budapest.
- Masako Watabe. 2024. A polylectal linguistic resource for Rromani. In Max Silberztein. (ed.) *Linguistic Resources for Natural Language Processing: On the Necessity of Using Linguistic Methods to Develop NLP Software*. pages 147-172. Springer, Cham.
- Max Silberztein. 2003-. NooJ manual. <https://nooj.univ-fcomte.fr>
- Max Silberztein. 2016. *Formalizing Natural Languages: the NooJ approach*. Wiley Ed.: Hoboken NJ.
- Medo Gurbetovski, Mozes Heinschink, and Daniel Krasa. 2010. *Guide de conversation rromani de poche*. ASSIMIL, Paris.
- Rajko Đurić. 2006. E rromani çhib. In Marcel Courthiade (ed.) *La littérature des Rroms, Sintés et Kalés*. pages 67-68. INALCO, Paris.
2010. *Atlas of the World's Languages in Danger*. UNESCO, Paris.
- Facebook. <https://www.facebook.com/>
- Google Translate. <https://translate.google.com/>
- NooJ platform. <https://nooj.univ-fcomte.fr/>
- ROMLEX. <http://romani.uni-graz.at/romlex/>
- Russian Romani Corpus. <http://web-corpora.net/RomaniCorpus/search/>
- La langue romani - un atout pour l'éducation et la diversité (exhibition). 2014. Council of Europe, Strasbourg.