

DEVELOPMENT OF TOOLS TO INFLECT COMPOUND WORDS IN ROMANIAN

Elena Boian¹, S. Cojocaru¹, A. Colesnicov¹, L.Malahova^{1*}, Tatiana Baltaga²

*Institute of Mathematics and Computer Science, Academy of Sciences of Moldova¹
State University of Moldova, Department of Mathematics and Information Technologies²*

e-mail: baltaga@usm.md, lena@math.md sveta@math.md kae@math.md

mal@math.md

ABSTRACT

One of basic functions at management of electronic dictionaries is its population with new words. The population of words can be manual when word by word is entered with morphologic attributes, or automated when a special tool is used to inflect simple and compound words. An approach to compound words inflexion in Romanian is discussed.

Keywords: morphologic dictionary, compound words, automated inflexion

INTRODUCTION

Morphologic dictionary contains words of a natural language with the attributes that are specific for this language morphology. Specific attributes make the essence of a morphologic dictionary. To generate the actual paradigm we use a program for automatic inflexion. The dictionary serves a primary open base to set new attributes for different domains of a natural language [1-2]. The dictionary structure for a highly inflective language is a difficult problem

One of basic functions at management of electronic dictionaries is its population with new words [3]. The population of words can be manual when word by word is entered with morphologic attributes (we can call it also the declarative mode), or automated when a special tool is used to inflect simple and compound words (procedural mode) [4–5]. The problem of the declarative mode is the necessity to have all word-forms ready for each new word. The procedural method can meet a problem if the inflexion rules for a new word were not programmed. The manual population is very difficult, because we need to enter word-forms and the corresponding morphological information. It is therefore naturally to try to automate the process.

The word inflexion program facilitates dictionary population. Information concerning part of speech (noun, verb, etc.) and maybe something else (e.g., gender for nouns) can be obtained

interactively. Then deep linguistic analysis with possible additional questions (usually about the alternation or suffixes) permits to predict the structure of the forms.

Special attention is given to nouns and adjectives declination, and verbs conjugation because they generate a lot of flexions. Pronouns, articles and numerals are not so numerous and may be entered in the dictionary manually. The problem of the complete formalization of the inflection process is very difficult, and we solved it partially. We ask the user for additional information in some special cases. This information concerns morphological categories, and consonant and vowel alternations.

A part that was not implemented till now is the automated inflexion of compound Romanian words.

Methods of word composition in Romanian are **prefixing, suffixing** and **colliding** [6–8]. Prefixing and suffixing form simple words. These words are inflected by the existing word inflexion program [4–5].

Another mode of word formation is **colliding**. The collided word units in itself constructive elements that were combined during language evolution, or two or more words, which produced a new word with a new contents [6-8].

The present word supposes as its purpose further development of existing tools of word inflexion for Romanian [4–5] and their generalization for compound words, and the use of the developed tools for population of a Romanian lexicon with word-forms and morphological information for compound words [1-2].

Word colliding

We mentioned above that word **colliding** is made in two modes. Table 1 shows the first mode, when a word is varied by additional auxiliary elements. The second mode is a colliding of two different words with their sense making a new word with a new content (Table 2).

Table 1. *The first mode of compound word formation*

Part of speech	Constructive elements	Examples
Common nouns	words that has in their structure elements, which can not be used independently: <i>aero-</i> , <i>balneo-</i> , <i>micro-</i> , <i>neuro-</i> , <i>-cid</i> , <i>-naut</i> , etc.	<i>aerodrom</i> , <i>balneolog</i> <i>microfon</i> , <i>neurology</i> , <i>genocid</i> , <i>astronaut</i> ,
	composed of two or more words where only a final component is inflected	<i>binecuvântare</i>
	composed of two or more words	<i>metaloplasie</i>
Adjective	composed of existing words or radices: adverb+ adjective <i>ne+ mai+ adjective</i>	<i>binefăcător</i> <i>nemaiauzit</i> , <i>nemaivăzut</i>
	composed of an auxiliary element and an existing word: <i>mono-</i> , <i>uni-</i> , <i>bi-</i> , <i>multi-</i> , <i>poli-</i> , <i>semi-</i>	<i>monocelular</i> , <i>monofazic</i> , <i>univalent</i> , <i>univoc</i> , <i>bianual</i> , <i>bisilabic</i> <i>multimi-lenar</i> , <i>multinațional</i> , <i>polisemantic</i> , <i>politehnic</i> , <i>semiovscur</i> , <i>semioficial</i>
Verb	composed of <i>bine+</i> verb	<i>binecuvînta</i>

	composed of <i>auto+</i> verb	<i>autoadministra</i>
	composed of <i>tele+</i> verb	<i>telecomanda</i>
	composed of adverb <i>nemai+</i> participle or gerund	<i>Nemaicunoscut, nemaipomenit, nemaiauzind, etc.</i>
Numeral	formed from two or more words	<i>noiăsprezece, treizeci</i>
Pronoun	words that has in their structure elements, which can not be used independently:	<i>dumneata, dumnealui, dumneavoastră</i> (polite personal pronouns); <i>sieși, sineși</i> (reflexive pronouns); <i>însumi, însuși, însuși, ș.a.</i> (identity pronouns); <i>aceiași, aceeași; ș.a.</i> (demonstrative pronouns); <i>careva, vreun, vreunul, vreuna, vreo</i> (indefinite pronouns)
Adverb	preposition + adverb:	<i>adeseori, apururea, deasupra, înapoi</i>
	preposition + cardinal numeral	<i>acasă, alături, alene, alocuri, deopotrivă, departe, deseară, deval, devreme, pesemne, etc.</i>
	:	
	prepositions <i>în + de +</i> another part of speech	<i>Intruna, delaolaltă, laolaltă; îndeaproape, îndeajuns, îndelung, înaemînă, îndeobște, îndeosebi, îndeșeară, etc.</i>
	adjective + noun	<i>deseori, rareori</i>
	pronoun + adverb	<i>astfel, altfel, alaltăieri, alaltăseară, altădată, altcîndva, altcum, altcumva; altcîndva, altcum, altcumva</i>
	adverb + adverb:	<i>bunăoară, etc.</i>
	adverb + auxiliary element	<i>cîndva, cumva, încotrova, undeva, oarecînd, oarecum, oarecît, oareunde, oricînd, oricum, oriîncotro, oriunde, oricît, etc.</i>
adverb + <i>și</i>	<i>cîtuși, iarăși, totuși</i>	
	several heterogeneous elements	<i>cîteodată, nicidecum, nicidecum, numaidecît, orișicînd, orișicum, orișunde, pasămite, încămite, etc.</i>
Preposition	composed from constructive elements combined in a single word	<i>despre, dinspre, dintre, deasupra, dindărătul</i>
Conjunction	composed from constructive elements combined in a single word	<i>așadar, deoarece, fiindcă, încît</i>

Table 2. The second mode of compound word formation

Part of speech	Constructive elements	Examples
Common noun	common noun formed from two or more independent common nouns, where the first element is inflected	<i>tonă-kilometru (tone-kitometru)</i>
	common noun formed from two or more independent common nouns, where the second element is inflected	<i>liber-cugetător (liberului-cugetător)</i>
	common noun formed from two or more independent common nouns, where all elements are inflected	<i>mașină-unealtă (mașinii-unelte)</i>
	common noun formed from two or more independent common nouns, where no elements are inflected	<i>apă-albă</i>
Proper noun	proper noun formed from two or more nouns in nominative-accusative case	<i>Grinăuți-Moldova</i>
	proper noun followed by another noun, being personal or geographic name	<i>Ana-Maria</i>
	proper noun followed or forwarded by a qualificating noun or adjective	<i>Alb-împărat</i>
	proper noun followed or forwarded by a common noun	<i>Vadul-Leca</i>

	other structure	<i>Lacul Sărat, Polul Nord</i>
Adjective	adjective + adjective	<i>alb-galbui,</i>
Numeral	formed of two or more words linked by prepositions, from which at least one of numerals is a simple one	<i>al doilea (a doua), patruzeci și opt cinci zecimi; de două ori, de trei ori, etc.</i>
Pronoun	with different elements of composition	<i>nici unul, nici un, nici una, nici o, etc.</i> (negative pronoun); <i>Alteța sa, Sanctitatea sa, etc.</i> (polite pronoun)
Adverb	noun + noun, noun + adverb, adverb + adverb, verb + verb, interjection + interjection	<i>calea-valea, an-vară, cîine-cîinește, vrînd-nevrînd, treacă-meargă, încet-încet, hodoronc-tronc, țac-pac, etc.</i>
	adjective <i>astă</i> + noun	<i>astă-iarnă, astă-toamnă, astă-vară, etc</i>
	<i>adverb + noun</i>	<i>azi-dimineață ieri-noapte, ieri-seară, etc.</i>
	preposition <i>după</i> + noun	<i>după-amiază, după-masă, după-prînz, etc.</i>
	<i>preposition</i> <i>întru, dintru + adverb</i>	<i>într-aceea, dintr-acolo, într-adevăr, într-adins, într-aiurea, etc.</i>
Preposition	composed from several simple prepositions	<i>de la, de către, de pe, de pe la, de după, de sub pe lângă, până pe la, până pe lângă, etc.</i>
	prepositional locations of different structure	<i>afară de, în afară, pe aproape de. de jur împrejurul</i>
Conjunction	compound conjunctions are many different structures	<i>fără ca, pentru ca, pentru că, pentru ce, de ce, ca și cum, ca și cînd, pe cîtă vreme, cît timp, ori de cîte ori, etc.</i>
Interjection	formed from one or more simple interjections	<i>cioc-cioc, hop-hop, hi-hi-hi, ha-ha-ha, tronca-tronca, etc.</i>
	repeating simple interjections	<i>miau,miau; mac,mac; boca,boca; etc.</i>

THE INFLEXION PROCESS

The inflexion process is as follows. The word to be inflected is entered and indicated as simple or compound. If the word is simple, the inflexion program works as described in [4–5]. If the compound word is formed in accordance with Table 1, the inflexion program works in the usual mode described in the [4–5]. If the compound word is formed as described in Table 2, additional information is needed to inflect this word. After inflection the manual correction of the resulting forms is possible. During manual correction the user has a possibility to edit the whole visible screen, however not inserting and deleting lines and without block operations. After the manual correction the words can be introduced into the dictionary.

CONCLUSIONS

The presented data and described methods are the first step to further development of automated tools for Romanian lexicon population, including the inflexion of compound words.

REFERENCES

1. Studii de Materiale privitoare la Formarea Cuvintelor in Limba Română Academia Republicii Populare Române. Institutul de Lingvistică din București. Volumul I. Editura Academiei Republicii Populare Române, București, 1959.
2. E. Boian, C. Ciubotaru, S. Cojocaru, A. Colesnicov, V. Demidova, L. Malahova: “Reusable linguistic resources for romanian.” - Second Conference of the Mathematical Society of the Republic of Moldova, August, 17-19, 2004, Chișinău, Republic of Moldova. – p. 54–57.
3. E. Boian, C. Ciubotaru, S. Cojocaru, A. Colesnicov, V. Demidova, L. Malahova: “Lexical resources for Romanian”. - în “Memoriile științifice ale Academiei Române”, București, România. 2004.
4. S. Cojocaru: Romanian Lexicon: Tools, Implementation, Usage. – In: Dan Tufis & Poul Andersen (eds.), Recent Advances in Romanian Language Technology. – ISBN 973–27–0626–0, Editura Academiei, 1997, I, pp. 107–114.
5. E. Boian, S. Cojocaru. The inflexion regularities for the Romanian language. – Computer Science Journal of Moldova, Vol. 4, Nr. 1 (10), 1995, pp.40–58.
6. E. Boian, S. Cojocaru, L. Malahova. Instruments pour applications linguistiques. La terminologie en Roumanie et en République de Moldova, Hors série N4, 2000.
7. Studii de Materiale privitoare la Formarea Cuvintelor in Limba Română Academia Republicii Populare Române. Institutul de Lingvistică din București. Volumul I. Editura Academiei Republicii Populare Române, București, 1959.
8. Gramatica limbii române.– București, Ed. Acad. Române, 1977.
9. Iorgu Iordan, Vladimir Robu. Limba Română Contemporană. Editura Didactică și Pedagogică, București, 1978.