

# RUSSIAN MORPHOLOGICAL ANALYSIS

**Serge A. Yablonsky**

Department of Computing, St.-Petersburg University of Transport

Russicon Company

Kazanskaya 56, Ap. 2

190000 St.-Petersburg, Russia

*e-mail: root@russicon.spb.su; tel/fax: +7-812-312-72-13*

## Abstract

In this paper the approach to the organization of Russian inflexion morphologic model and its application for the Russian language morphological analysis and disambiguation are described. We are concerned with the pos tagging of 150-million-word Russian corpora. The approach is particularly dependent on the language processor Russicon, and on wide usage of Russicon's electronic dictionaries.

## 1. Introduction

Up-to-date language technologies contain efficient morphological analyzers for Romance, Germanic (Karttunen, 1983; Karttunen, Koskeniemi, Kaplan, 1987; Varile, Zampolli, 1996; Zaenen, Uszkoreit, 1996) and some Slavic (Chanod, 1997) languages. In the last 15 years Russian computational morphology has advanced at a great rate from first quite restricted systems towards large-scale practical morphological analyzers (Ashmanov I., 1995; Belonogov, Zelenkov, 1989; Belyaev, Surcis, Yablonsky, 1993; Bolshakov, 1990; Mikheev, Liubushkina, 1995; Segalovich, 1995). Now Russian word-form morphological analyzers are able to:

- classify an input string as a valid word of the supported lexicon and categorize it morpho-syntactically: part-of-speech (pos) category, number, gender, etc.

- generate a form of a word in accordance with certain morpho-syntactic features.

Although all these systems have an impressive lexicon (more than 1,500,000 word-tokens and 80,000 – 200,000 word paradigms) and are computationally efficient, few of them have user friendly interface to extend their lexicon with new words (Mikheev, Liubushkina, 1995; Yablonsky, 1990, 1998) and robust analysis of unrestricted Russian texts.

Russian part-of-speech (pos) tagging provides developing of disambiguation algorithms on the top of morphological text analysis. Alternate morphological analyses occur because of high categorial homonymy of inflective language. Only the sentential context used by disambiguation normally decides which analysis is appropriate (Zaenen, Uszkoreit, 1996).

In this paper the approach to the organization of Russian inflexion morphologic model and its application for the Russian language morphological analysis and disambiguation are described. Our purpose is to mark each word of 150-million-word Russian corpora (Yablonsky, 1998, a; 1999) with a pos and lexical categories tags.

The approach is particularly dependent on the language processor Russicon, and on wide usage of Russicon's electronic dictionaries (Belyaev, Surcis, Yablonsky, 1993; Yablonsky, 1990, 1997, 1998).

## 2. General model of inflection morphology

For the formal description of *inflection morphology model* the set theory is used. It is one of the best ways for description of Russian inflection morphology (Bider, Bolshakov, 1976; Kulagina, 1986).

Here we present the general set model that permits to define mostly all sides of inflection morphology. It is realized in the morphological analyzer of the Russicon language processor (Yablonsky, 1990; Belyaev, Surcis, Yablonsky, 1993; Yablonsky, 1998). In this model some concepts are introduced for the first time and other have new or more full meaning.

### 2.1. Definitions

Lexicon of the inflective language is named  $Q$   $= \{q_1, q_2, \dots, q_{N_q}\}$ , where  $q_i, i = \overline{1, N_q}$  — legal sequence of alphabet  $L$  characters named as a *word*. Set  $Q$  is an infinite denumerable set, because lexicon is permanently enlarging. Denurable system  $W$  of in general intersecting finite subsets  $W_i$  named as lexemes is given in  $Q$ ,  $\bigcup W_i = Q$ . In the unstrict, lexeme  $W_i = \{w_{i,1}, w_{i,2}, \dots, w_{i,N_{wi}}\}$ , where  $w_{i,j} \in Q$ ,  $j = \overline{1, N_{wi}}$ , is a set of words that has the same interpretation and etymology. Words  $w_{i,j} \in W_i$  are named as *word forms or inflective forms of lexeme  $W_i$* .

Let us introduce for each inflection form  $w$  of lexeme  $W_i$  number  $N_{wis}$  from the interval  $1 \div |w|$ . The combination of a surface form and its analysis as a canonical form and inflection is called a *lemma*. *Paradigm* of the word is composed of all inflective forms of one lemma.

For every word one or several paradigms could be mapped:

$$W_i \rightarrow \overline{\beta}_m (\beta_{m,1}, \beta_{m,2}, \dots, \beta_{m,p}).$$

The number of word forms  $P$  for inflective

language is rather high ( $P \gg 1$ ) for some parts of speech. For example, in Russian language  $P > 100$  for verbs.

The part of the word including prefix(es), root and suffix(es) is called *word inflective stem (WIS)*. The length  $L_{wis}$  of WIS is calculated  $L_{wis} = |WIS| = L_{word} - L_{end}$ , where  $L_{word}$  is the length of the word and  $L_{end}$  — the length of the ending.

For the character sequences  $WIS = (b_1, \dots, b_{N_{wis}})$ , and  $INFL = (b_{N_{wis}+1}, \dots, b_{N_b})$  — *ending of inflection form (inflection)*, condition  $w = WIS \oplus INFL$  is fulfilled. Operation of two lists concatenation is denoted by ' $\oplus$ ';  $|WIS| = L_{wis} = 1 \div |w|$ ;  $|INFL| = 0 \div (|w| - 1)$ . Null length inflection is called *empty* and is denoted by '+', and corresponding flexion is denoted by '-'. WIS of the lemma is denoted by  $WIS^*$ .

The part of the word including prefix(es) and root is called *word formative stem (WFS)*. The length  $L_{wfs}$  of WFS is calculated  $L_{wfs} = |WFS| = L_{word} - L_{suf} - L_{end} = L_{wis} - L_{suf} = L_{pref} - L_{root}$ , where  $L_{suf}$  is the length of the suffix(es),  $L_{pref}$  — the length of the prefix(es),  $L_{root}$  — the length of the root(s). The range of the introduced parameters are:  $L_{root} = 1 \div L_{wis}$ ,  $L_{pref} = 0 \div (L_{wis} - 1)$ ,  $L_{suf} = 0 \div (L_{wis} - 1)$ ,  $L_{wfs} = 1 \div L_{wis}$ .

The morphological analyzer has two main parts:

- dictionary with linguistic knowledge of the language;
- program realization of morphologic model's algorithms.

In general dictionary is represented in such form:

$$V_i = \{W_i, f_i\}, \quad i = \overline{1, N_v},$$

where  $W_i = (a_1, a_2, \dots, a_{L_i})$  — lexical part of dictionary's article: the word or phrase, composed from the alphabet characters  $A = \{a_s : s = 1, \dots, N_a\}$ ; tag part  $f_i = (f_1, f_2, \dots, f_k)$  — subset of tags from the set  $F = \{f_r : r = 1, \dots, N_f\}$ ,  $N_v$  — number of the words (word-tokens)

in the dictionary, for large-scale dictionaries of inflective languages usually  $N_v > 1500\ 000$ .

## 2.2. Morphological analysis model

Let  $H = \{h_1, h_2, \dots, h_{N_h}\}$  be the set of part-of-speech (pos) categories and  $P = \{p_1, p_2, \dots, p_{N_p}\}$  — lexical categories (LC) of gender, number etc. Each element  $p_i \in P$ , where  $i = \overline{1, N_p}$ , represents the set of concrete realizations of lexical category  $p_i = \{p_{i,1}, p_{i,2}, \dots, p_{i,N_i}\}$ . Let us choose one element in  $P$  (for definiteness  $p_1$ ) named *type* and denoted by  $T$  ( $T = p_1, T \in P$ ),  $T = \{t_1, t_2, \dots, t_{N_t}\}$ .

For example, Russian language model includes:  $H^1 = \{h_1 = \text{"noun"}, h_2 = \text{"adjective"}, h_3^2 = \text{"verb"}, h_4 = \text{"particle"}, h_5 = \text{"parenthetic word"}, h_6 = \text{"modal word"}, h_7 = \text{"adverb"}, h_8 = \text{"conjunction"}, h_9 = \text{"interjection"}, h_{10} = \text{"preposition"}, h_{11} = \text{"abbreviation"}, h_{12} = \text{"unit of measure"}, h_{13} = \text{"pronoun"}, h_{14} = \text{"numeral"}, h_{15} = \text{"adverbial participle"}, h_{16} = \text{"composition or special prefix"}\}$ .  $P = \{p_1 = \text{"case"}, p_2 = \text{"gender"}, p_3 = \text{"number"}, p_4 = \text{"time"}, p_5 = \text{"person"}, p_6 = \text{"degree"}, p_7 = \text{"voice"}, p_8 = \text{"aspect"}, p_9 = \text{"mood"}, p_{10} = \text{"form"}, p_{11} = \text{"transitivity"}, p_{12} = \text{"reflexive"}, p_{13} = \text{"animate"}\}$ , where  $p_1 = \{\text{"nominative"}, \text{"genitive"}, \text{"dative"}, \text{"accusative"}, \text{"instrumental"}, \text{"prepositional"}\}$ ,  $p_2 = \{\text{"masculine"}, \text{"feminine"}, \text{"neuter"}, \text{"masculine/feminine"}\}$ ,  $p_3 = \{\text{"singular"}, \text{"plural"}\}$ .  $p_4 = \{\text{"present"}, \text{"past"}, \text{"future"}, \text{"present / future"}\}$ ,  $p_5 = \{\text{"1st person"}, \text{"2nd person"}, \text{"3rd person"}\}$ ,  $p_6 = \{\text{"superlative"}, \text{"comparative"}\}$ ,  $p_7 = \{\text{"active"}, \text{"passive"}\}$ ,  $p_8 = \{\text{"imperfective"}, \text{"perfective"}, \text{"perfective and imperfective"}\}$ ,  $p_9 = \{\text{"indicative"}, \text{"imperative"}\}$ ,  $p_{10} = \{\text{"full"}, \text{"short (predicative)"}, \text{"infinitive"}\}$ ,  $p_{11} = \{\text{"transitive"}, \text{"intransitive"}\}$ ,  $p_{12} = \{\text{"reflexive"},$

$\text{"irrevocable"}\}$ ,  $p_{13} = \{\text{"animate"}, \text{"inanimate"}\}$ .

We take that  $(\forall h_k : h_k \in H) (\exists T_k : T_k \subset T, T_k \neq \emptyset)$ , i.e. at least one type exists for each part of speech and is named *ordinary*, and also  $T_k = \{t_{k,1}, t_{k,2}, \dots, t_{k,N_k}\}$ , where  $k = \overline{1, N_h}$ .

For example, in Russian for  $h_1 = \text{"noun"}$ ,  $T_1 = \{\text{"ordinary"}, \text{"invariable"}, \text{"substantival"}\}$ . For  $(h_k, t_{k,j}) = (\text{"noun"}, \text{"ordinary"})$ :  $P_{1k,t} = \{\text{"gender"}\}$ ,  $P_{2k,t} = \{\text{"case"}, \text{"number"}\}$ ,  $P_{3k,t} = \{\text{"animate"}\}$ ,  $X_{k,t} = (X^*_{k,t} = \{\text{"nominative case"}, \text{"singular number"}, \{\text{"genitive case"}, \text{"singular number"}, \{\text{"dative case"}, \text{"singular number"}, \{\text{"accusative case"}, \text{"singular number"}, \{\text{"instrumental case"}, \text{"singular number"}, \{\text{"prepositional case"}, \text{"singular number"}, \{\text{"nominative case"}, \text{"plural number"}, \{\text{"genitive case"}, \text{"plural number"}, \{\text{"dative case"}, \text{"plural number"}, \{\text{"accusative case"}, \text{"plural number"}, \{\text{"instrumental case"}, \text{"plural number"}, \{\text{"prepositional case"}, \text{"plural number"}, \{\text{"2-nd genitive case"}, \text{"singular number"}, \{\text{"2-nd instrumental case"}, \text{"singular number"}, \{\text{"2-nd prepositional case"}, \text{"singular number"}, \{\text{"2-nd accusative case"}, \text{"singular number"}, \{\text{"2-nd accusative case"}, \text{"plural number"}\}\})$ .

Then

$$(\forall h_k \forall t_{k,j} : h_k \in H, t_{k,j} \in T_k, k = \overline{1, N_h}, j = \overline{1, N_k})$$

$$(\exists P_{k,t} : P_{k,t} \subset P, \bigcup_{k=1}^{N_h} \bigcup_{t=1}^{N_k} P_{k,t} = P),$$

i.e. for each part of speech exists its own, may be empty, set of LC. Elements of  $P_{k,t}$  are named LC of  $t$ -type  $h_k$ . For all  $P_{k,t}$  exists partition on three nonoverlapping and may be empty subsets, named  $P^1_{k,t}, P^2_{k,t}, P^3_{k,t}$ :

$$(\forall P_{k,t} : P_{k,t} \subset P) (\exists P^1_{k,t} \exists P^2_{k,t} \exists P^3_{k,t} : P^1_{k,t} \cup P^2_{k,t} \cup P^3_{k,t} = P_{k,t}, P^1_{k,t} \cap P^2_{k,t} \cap P^3_{k,t} = \emptyset).$$

Elements of  $P^1_{k,t}$  are named as *ordinary LC*, elements of  $P^2_{k,t}$  — *special LC*, elements of  $P^3_{k,t}$  — *individual LC* of  $t$ -type  $h_k$ . For each  $P_{k,t}$  set, if  $P_{k,t} \neq \emptyset$ , there exists ordered sequence of sets  $X_{k,t} = (X^1_{k,t}, X^2_{k,t}, \dots, X^{N_{k,t}}_{k,t})$ . That is, if

<sup>1</sup> Our model for Russian slightly differs from the classic: we include in the sets  $H$  and  $P$  some additional elements.

<sup>2</sup> In the paradigm of the verb we include participle and adverbial participle.

$P_{k,t}^2 = \{p_{k,t,1}^2, p_{k,t,2}^2, \dots, p_{k,t,N_{kt}^2}^2\}$ , then  $X_{k,t}^1 = \{x_{i,j} : x_{i,j} \in p_{k,t,i}^2, i = \overline{1, N_{kt}^1}\}$ , where  $1 = \overline{1, N_{kt}^1}$ , and if  $P_{k,t}^2 = \emptyset$ , then it is considered that  $X_{k,t} = (\emptyset)$ .

We shall call sequence  $X_{k,t}$  the *sequence of lexical categories of word inflective paradigm of t-type*  $h_k$ . One of lexical categories, usually the first, is named  $X_{t,k}^*$  and called *normalized*.

There exists a single pair  $(h_k, t_{k,j})$ ,  $(h_k \in H, t_{k,j} \in T_k, k = \overline{1, N_h}, j = \overline{1, N_k})$  for every lexeme and, therefore, ordered list of LC  $X_{k,t}$ .

Let us define function  $f_{l \rightarrow W}(l)$ ,  $l = \overline{1, N_X}$ , with range of values  $W_l = W_l \cup \{\varepsilon\}$ , where  $\varepsilon$  – *dummy*, nonexistent word form. Thereby, for every lexeme  $W_l$  an ordered sequence of word-forms  $Y_{W_l} = (y_1, y_2, \dots, y_{N_X})$  ( $y_j \in W_l$  for  $j = \overline{1, N_X}$ ) could be formed. Such sequence is called *word changing paradigm of lexeme*  $W_l$  (*WCP*). If for lexeme  $W_l$  exists such  $l$ , that  $f_{l \rightarrow W}(l) = \varepsilon$ , it is said that lexeme  $W_l$  has a *dummy word changing paradigm* (Apresyan, 1989).

If the pair  $(h_k, t_{k,j})$  corresponds to lexeme  $W_l$  and for some  $l = l^*$  from  $\overline{1, N_X}$  conditions:  $f_{l \rightarrow X}(l^*) = X_{k,t}^*$  and  $y^* = f_{l \rightarrow W}(l^*)$  ( $y^* \in Y_{W_l}, y^* \neq \varepsilon$ ), are fulfilled, then we shall call the word form  $y^*$  as *normalized form or lemma of lexeme*  $W_l$ . Usually  $y^* = y_1$ . As a rule, infinitive is a lemma for the verb etc.

Let  $Y_{W_l}$  — WCP of lexeme  $W_l$ . Then word form's inflections of paradigm  $Y_{W_l}$  form ordered sequence denoted by  $Y_{FLCi}$ . Inflection class (FC) number  $I$  denoted by  $FC_I$  is the five:

$$FC_I = \langle h_k, t_{k,j}, P_{k,t}^1, X_{k,t}, Y_{FLCi} \rangle,$$

where  $h_k$  – some part of speech;  $t_{k,j}$  – some realization of LC *type* for corresponding part of speech;  $P_{k,t}^1$  – ordinary LC, corresponding to  $t_{k,j}$ ;  $X_{k,t}$  – sequence of special LC of WCP, corresponding to  $t_{k,j}$ ;  $Y_{FLCi}$  – some  $I$ -th sequence of inflections, also called WCP of

FC, where  $|X_{k,t}| = |Y_{FLCi}|$ . Inflection class concept was first used by (Belonogov, Zelenkov, 1985), although inflection class was understood only as ordered sequence of inflections.

Let for lexeme  $W$   $WIS^* = (b_{*1}^*, b_{*2}^*, \dots, b_{*N_{WIS}^*}^*)$  and exists  $WIS_m = (b_1, b_2, \dots, b_{N_{WIS}})$ , where  $m = 1 \div |X_{k,t}|$ , such, that  $WIS_m \neq WIS^*$ . Consequently, exists natural number  $N_0$ ,  $N_0 = 0 \div \min(N_{WIS}, N_{WIS}^*)$ , such, that  $(b_{*1}^*, b_{*2}^*, \dots, b_{*N_0}^*) = (b_1, b_2, \dots, b_{N_0})$ ;  $(b_{*N_0+1}^*, \dots, b_{*N_{WIS}^*}^*) \neq (b_{N_0+1}, \dots, b_{N_{WIS}})$ . Let us call the ordered sequence  $z_{l,m}^s = ((b_{N_0+1}, \dots, b_{N_{WIS}}), (b_{*N_0+1}^*, \dots, b_{*N_{WIS}^*}^*))$ , allowing to obtain lemma  $WIS$  from some word form  $WIS$ , *direct substitution*. Here  $I$  is a FC number,  $m$  – position number in  $WIS$  FC,  $s$  – exact pair number among other pairs in the  $m$ -th position. For each  $I$ -th FC is defined ordered set  $Z_I$  (may be empty):  $Z_I = \{z_{l,1}, z_{l,2}, \dots, z_{l,N_{zi}}\}$ , где  $N_{zi} = |X_{k,t}|$ . Each  $z_{l,m} = \{z_{l,m}^1, z_{l,m}^2, \dots, z_{l,m}^{N_{zim}}\}$ , where  $m = 1 \div |X_{k,t}|$ , also is a set of pairs of direct substitutions (may be empty). If the pair  $z_{l,m}^s = (b^m, b^*)$  is a direct substitution, then the pair  $(b^*, b^m)$  is called *reverse substitution*. Reverse substitution allows obtaining some  $m$ -th word form  $WIS$  from lemma  $WIS$ . There is one-to one correspondence between the sets  $B^* = \{\dots, b^*, \dots\}$  and  $B^m = \{\dots, b^m, \dots\}$ . Thus,  $|B^*| = |B^m|$ ; if  $(b^*, b^{m_1})$  and  $(b^*, b^{m_2})$ , then  $b^{m_1} = b^{m_2}$ ; if  $(b^{*1}, b^m)$  and  $(b^m, b^{*2})$ , then  $b^{*1} = b^{*2}$ . The letters from the constant part of  $WIS$  could be added to the beginnings of such character sequences for achievement of this term.

For example, the genitive of the plural noun *КОПЕЙКА* (copeck) with lexeme  $WIS^* = (КОПЕЙК)$  is  $WIS_7 = (КОПЕЕК)$ . Direct substitution should be (ЕК, ЙК), but for the lexeme of the same inflexion class (FC = 154) «ПУЛЬКА» (kitty or pellet or pool) direct substitution in the same position must be (ЕК, ЪК). This generates ambiguity. Therefore, two pairs of direct substitutions: (ЕЕК, ЕЙК) и (ЛЕЕК, ЛЬК) are formed in the

morphology model for inflexion class 154 and  $m = 7$ . Thus, for some inflection classes the set of direct substitutions should be formed.

So, for obtaining word form of WCP with given LC it is enough to define WIS of the lemma, number of the inflexion class and the number of word form in WCP, thus the three  $\langle \text{WIS}^*, \text{FC}, \text{I} \rangle$ . If  $Y_1 = \text{'-'}'$ , then for given FC and, accordingly, for given lexeme the word form with such LC does not exist. However, even if  $Y_1 \neq \text{'-'}'$ , paradigm of the given lexeme could be dummy. Such situation is described with the help of the set  $P^3_{k,t}$  of individual LC of given lexeme.

For example, lexemes «ДЕЛАТЬ» (do) and «СДЕЛАТЬ» have the same inflection class 175 and, hence, the same realization of ordinary and special LC, but they have different value of aspect: verb «ДЕЛАТЬ» – imperfective aspect, verb «СДЕЛАТЬ» – perfective aspect. So LC aspect should be the individual LC for this pair. Additionally, the individual LC could impose restriction on the existence of some inflections of the word. In the above example for  $\text{FC} = 175 \text{ FLC}_{44} = \langle \text{ЫЙ} \rangle$  и  $Z_{175,44} = \langle \text{ЕМ}, \sim \rangle$ , where sign  $\sim$  designates empty sequence. For the verb «ДЕЛАТЬ»  $\text{WIS}^* = \langle \text{ДЕЛА} \rangle \rightarrow y_{44} = \langle \text{ДЕЛАЕМЫЙ} \rangle$ . For the verb «СДЕЛАТЬ»:  $\text{WIS}^* = \langle \text{СДЕЛА} \rangle \rightarrow y_{44} = \langle \text{СДЕЛАЕМЫЙ} \rangle$ . This contradicts with Russian language standard.

So in the morphologic model should be the rules “rejecting” some inflection forms according their individual LC information. Such exclusion for given lexeme could be set explicitly by indicating the number of concrete inflection.

For example, for lexeme «МЕЧТА» (dream) there is no  $y_8$  – plural genitive inflection. The set of individual LC realizations of lexeme inflections and numbers of forbidden inflections of WIP are considered to be individual feature of lexeme and are marked by I.

Thus, LC of every lexeme  $W_i$  could be given by three:

$$W_i = \langle \text{WIS}^*_i, \text{FC}, \text{I}_i \rangle \quad (1).$$

### 2.3. Derivational morphology and compounding

Derivational morphology is based on detection of fixed expressions (more than 2000 of Russian idioms, proverbs, sayings), multiword prepositions, prefixes/suffixes with strong derivation functions and productive central derived forms (Kuznetcova, Efremova, 1986; Efremova, 1996), compounds (3000 of most frequent Russian compounds), processing 198 features consisting of morpho-syntactic features, derivational features, stylistic features and punctuator features.

## 3. Russian morphological dictionary

General lexicon of the dictionary is formed from the intersection of such lists (Yablonsky, 1998):

- Russian basic grammatical dictionary (80 000 word paradigms);
- Russian thesaurus (8 696 synonym rows, word list containing approximately 30 000 word paradigms);
- Large Russian explanatory dictionary<sup>3</sup> (more than 130 000 word paradigms from the language of the Eighties and the beginning of the Nineties);
- Orthographic dictionary<sup>4</sup> (60 000 word paradigms);
- Computer dictionary (1 500 word paradigms);
- Geographical names dictionary (1 500 word paradigms);

<sup>3</sup> Russicon company has the copyright on electronic version of the paper dictionary: Kuznetcov S.A.

(Large Explanatory Dictionary). St.-Petersburg. 1998, 1536 p. (in Russian)

<sup>4</sup> Russicon company licensed electronic version of the paper dictionary: Solov'ev N. V. Orthographic dictionary plus orthographic and punctuation reference guide, S.-Petersburg, 1996.

- Russian personal names, patronymics and surnames dictionary (10 000 word paradigms of Russian personal names, diminutives and patronymics, surnames of the world famous and Russian famous people);
- Business dictionary (1 500 word paradigms);
- Juridical dictionary (1 500 word paradigms);
- Jargon dictionary (5 000 entry word paradigms).

Joint number of all different Russian word paradigms is approximately **160 000**. Comprehensive description of all mentioned dictionaries is done in (Yablonsky, 1998).

## 3.2. Compressed database of Russian morphological dictionary

### 3.2.1. WFS - dictionary

In the compressed WFS - dictionary database all WIS are distributed into word forming groups (WFG). Word forming group consists of such set of fours:

$$\langle \text{WFS}_i, \text{SUF}, \text{FC}, \text{I}_i \rangle (2),$$

where SUF – suffix (number of the suffix), FC – inflection class number;  $\text{WIS}_i^* = \text{WFS}_i \oplus \text{SUF}$ . Only first 255 maximum frequent suffixes are coded as separate linguistic units in compress WFS-dictionary realization. Other suffixes are included in WFS. Thus WIS are distributed into **42 874** WFS. Capacity of the compressed dictionary is **990 K**.

### 3.2.2. WIS -realization

For increasing speed of morphological analysis all WIS with stem gradation were generated. This formed **179 289** word inflection stems for database. In the compressed WIS - dictionary database the ordered sequence of all lexemes represented by (1) is stored. The speed of analysis is increased in 10 times.

Besides, we use several additional tables: table of inflection classes, inflection class — inflections, inflection — inflection classes,

inflection class — right direct substitutions, joint right direct and right inverse substitutions, direct and inverse tables of suffixes, prefixes and substitutions in prefixes, and some other.

## 4. Evaluation of the morphological analyzer

Morphologic analyzer has been tested on various texts including more than 50 million words of literary and newspaper texts, Russian laws (1990–1995 years) available from Russicon text corpora (Yablonsky S.A., 1999, a, b). The results demonstrated right recognition of 95 – 98 % of text words. For Russian language the morphological analyzer leaves 1-5 % of all words in running text without the correct analysis when 10-15 % of words still have two or more analyses because of high categorial homonymy of inflective language. Only the sentential context used by disambiguation normally decides which analysis is appropriate.

C-realizations of morphologic analyzer and disambiguater are currently available for MS DOS and Windows 9X/NT.

Detailed description of Russicon Slavonic language processor could be find in forthcoming Belyaev B.M., Surcis A.S., Yablonsky S.A. / Yablonsky S.A. (ed.), (1999). Slavonic Language Processor RUSSICON, St.-Petersburg, Russia.

## 5. Concluding remarks

Now we are developing a rule-based disambiguater that uses Russian constraint grammar that consists of 1000 disambiguation rules, syntactic markers of lexemes.

Syntactic markers consist of government, concord and parataxis models for verbs, nouns, adjectives, prepositions and other parts of speech (Apresyan, 1985, 1989; Crockett, 1975; Iomdin, 1990; Russian grammar, 1980). To

further improve of the disambiguating power of the tagger, the grammar and the number of syntactically marked lexemes is to be extended.

Russian disambiguator is now tested against several kinds of Russian texts from Russicon collection (Yablonsky, 1999, a, b). The first results (500 000 words) showed that tagger has fully disambiguated (no homonymy) from 76% (contemporary Russian artistic avant-garde texts) up to 90 – 95% “neutral” Russian texts.

## 5. Acknowledgements

My thanks go to university and Russicon company colleagues Boris Belyaev, Anatol Surcis and Michael Kazakov. We worked together on a common design.

## 6. References

- Apresyan Y.D., (1985) Синтаксические признаки лексем (Syntactic markers of lexeme). *Russian Linguist.* Vol. 9, №2/3, p. 289–318. (in Russian)
- Apresyan Y.D., ed. (1989) Лингвистическое обеспечение системы ЭТАП-2 (Linguistic framework of the system ЭТАП-2). Moscow. (in Russian)
- Ashmanov I. (1995) Grammar and Style Checker for Russian Texts. In *Proceedings of Dialog'95 International Workshop on Computational Linguistics and its Applications*. Kazan, Russia.
- Belonogov G.G., Kuznezov B.A. (1983) Языковые средства информационных систем (Language tools for information systems). Moscow, 288 p. (in Russian)
- Belonogov G.G., Zelenkov Y.G. (1985) Алгоритм морфологического анализа русских слов (An Algorithm for Morphological Analysis of Russian words). In journal “Issues of information theory and practice”, №53. Moscow. (in Russian)
- Belyaev B.M., Surcis A.S., Yablonsky S.A. (1993) Russian Language Processor RUSSICON: Design and Applications. In the *Proceedings of the East-West Artificial Intelligence Conference EWAIC-93*, Moscow, pp.175-180.
- Bider I.G., Bolshakov I.A. (1976) Формализация морфологического компонента модели “Смысл $\leftrightarrow$ Текст”. I. Постановка проблемы и основные понятия (Formalization of morphological component within the Meaning $\leftrightarrow$ Text framework). *Reports of USSR Academy of Science on Technical Cybernetics*. №6, pp.42–57. (in Russian)
- Bolshakov I.A. (1990) A Large Russian Morphological Vocabulary for IBM Compatibles and Methods of its Compression. In the *Proceedings of the 13th International Conference on Computational Linguistics COLING-90*. Helsinki, Finland.
- Belyaev B.M., Surcis A.S., Yablonsky S.A. / Yablonsky S.A. (eds). (1998). Forthcoming *Slavonic Language Processor RUSSICON*. St.- Petersburg.
- Chanod J. (1997) Current development for Central and Eastern European Languages. In *Proceedings of the Second European Seminar “Language Applications for a Multilingual Europe”*, Mannheim/Kaunas.
- Crockett D.B. (1975) Agreement in contemporary standard Russian. Cambridge (Mass.). 458 p.
- Denisov P.N., ed. (1978) Учебный словарь сочетаемости слов русского языка (Russian word combinability dictionary). Moscow. (in Russian)
- Efremova T.F. (1996) Толковый словарь слово-образовательных единиц русского языка (Explanatory dictionary of Russian language word forming units). Moscow. 636 p. (in Russian)
- Iomdin L.L. (1990) Автоматическая обработка текста на естественном языке: модель согласования (Automatic Natural text Processing: a Model of Grammatical Agreement). Moscow. 168 p. (in Russian)
- Kuznetcova A.I., Efremova T.F. (1986) Словарь морфем русского языка (Russian morpheme dictionary). Moscow. (in Russian)
- Karttunen L. (1983) KIMMO: a general morphological processor. In Dalrymple et al (Eds.). *Texas Linguistic Forum*, 22, Department of Linguistics, University of Texas at Austin, pp166-186
- Karttunen L., Koskeniemi K., Kaplan R. (1987) A Compiler for Two-Level Phonological Rules. Technical Report. Center for the Study of Language and Information. Stanford University.
- Kulagina O.S. (1986) Морфологический анализ русских именных словоформ (Morphologic analysis of Russian word forms). Internal Publication of the IPM. Moscow, Academy of Sciences, №10, 26p. (in Russian)
- Kuznetsov S.A. (1998) Большой толковый словарь русского языка (Large Explanatory Dictionary). St.-Petersburg. 1536 p. (in Russian)

- Mikheev A.S., Liubushkina L.A. (1995) Russian Morphology: An Engineering Approach Natural Language Engineering 1 (3), Cambridge University Press, pp. 235–263.
- Varile, G. B., Zampolli, A. (1996) *Survey of the State of Art in Human Language Technology*. Cambridge: Cambridge University Press.
- Popov E. V. (1986) Talking with Computers in Natural Language. Springer-Verlag, 305p.
- Russian grammar (1980). Moscow. Vol. 1,2. (in Russian)
- Rybakov F.I., Rudnev E.A., Petukhov V.A. (1980) Автоматическое индексирование на естественном языке. (Automatic indexing using natural language). Moscow. (in Russian)
- Segalovich I.S. (1995) Indexing of Large Russian Texts with a Dictionary Built Around the Sparse Hash Table. In Proceedings of Dialog'95 International Workshop on Computational Linguistics and its Applications. Kazan, Russia.
- Tapanainen P., Jarvinen T. (1994) Syntactic Analysis of Natural Language Using Linguistic Rules and Corpus-based Patterns. In *Proceedings of Fifteenth International Conference on Computational Linguistics (COLING '94)*. Kyoto, Japan.
- Tikhonov A.N. (1985) Русский словообразовательный словарь (Russian word forming dictionary). Moscow. Vol.1,2. (in Russian)
- Yablonsky S.A. (1990) Russian Language Processor RUSSICON. In *Actual problems of computer linguistics*, Tartu, Estonia. (in Russian)
- Yablonsky S.A. (1997) New Capabilities for Russian and Ukrainian Language Learning Based on the Language Processor Russicon. In Sake Jager, John Nerbonne and Arthur van Essen (eds). Forthcoming “*Language Teaching and Language Technology*”. Lisse: Swets and Zeitlinger.
- Yablonsky S.A., (1998). Russicon Slavonic Language Resources and Software. In: W. Teubert, E. Tognini Bonelli and N. Volz (eds.) *Proceedings of the Third European Seminar Translation Equivalence*, ( pp. 217-227), Montecatini Terme, Italy.
- Yablonsky S.A. (1999, a). Russian Written Language Corpora Development. In: Proceedings of the International Seminar Dialog99, May 30-June 8, Tarussa, Russia.
- Yablonsky S.A. (1999, b). Russian 20th Century Literature Digital Library for Language Teaching. In: Proceedings of the International Digital Libraries for Humanities Scholarship and Teaching JUNE 9-13, 1999, (pp. 252 - 253), University of Virginia Charlottesville, Virginia, USA.
- Zaenen A., Uszkoreit H. (1996) Language Analysis and Understanding. In *Survey of the State of the Art in Human Language Technology* (<http://www.cse.ogi.edu/CSLU/HLTsurvey/>).