

The First One-Million Corpus for the Belarusian NooJ Module

Ivan Reentovich¹(✉), Yuras Hetsevich¹, Valery Voronovich², Evgenia Kachan²,
Hanna Kozlovskaya², Angelina Tretyak², and Uladzimir Koshchanka³

¹ United Institute of Informatics Problems, Minsk, Belarus
mwshrewd@gmail.com, Yury.Hetsevich@gmail.com

² Belarusian State University, Minsk, Belarus
gamrat.vvv@gmail.com, evgeniakacan@gmail.com,
malavita3000@gmail.com, angelina_tret@mail.ru

³ The Centre for the Belarusian Culture, Language and Literature Research,
National Academy of Sciences of Belarus, Minsk, Belarus
koshul@gmail.com

Abstract. In this article the first one-million corpus for the Belarusian NooJ module is represented. The given corpus has been built up of texts, patched up into sections by different subject categories. From the broad list of possible subject categories in the sections the corpus focuses on fiction, historic, medical, scientific, sociological literature, etc. Given a great number of similar subject categories, the first one-million corpus can be considered as a first subject collection of texts for the Belarusian NooJ module.

The text corpus is expected to be suitable for research in the following aspects: word polysemy processing of various texts, polysemic punctuation marks processing, and a new lexical items search.

The first one-million corpus for the Belarusian NooJ module can be fully applicable in many fields of linguistic research.

Keywords: Corpora · Belarusian NooJ-module · Statistical analysis · Part-of-Speech tagging · Machine-learning algorithms · Levenshtein algorithm · The machine-learning model · Counter-check · Spelling errors · Concatenation-in-paradigm · Unknown words search · Known words search · Clustering · Belarusian N-Corpus · Text processing

1 Introduction

The chosen subject “First One-million Corpus for the Belarusian NooJ module” is among the most important, integral parts of future research in the field of speech recognition and synthesis. It is really a great step on the way to new investigations of the Belarusian language in the world of NooJ and linguistics.

The purpose of this paper is to introduce the elaboration, creation stages as well as stages of “deep” analysis and practical application of the first one-million corpus for the Belarusian NooJ module in the context of different aspects and approaches.

Besides, the first one-million Belarusian corpus for the Belarusian NooJ module will be applicable for solving the following fundamental tasks: optimizing and expanding the development of high-quality linguistic algorithms for the electronic text pre-processing in the TTS (Text-to-Speech) system.

Two Belarusian corpora were developed for NooJ [1] – the 1-VERSION corpus (1 million corpus.noc) and the MAIN corpus (First 1MLN Corpus for the Belarusian NooJ Module.noc). To make the process of corpus creation more productive, a special (descriptive) algorithm has been worked out.

2 Descriptive Algorithm for the First One-Million Corpus for the Belarusian NooJ Module

The main work on corpus compilation and analysis with the help of this algorithm was fulfilled on the basis of 1-VERSION corpus (Table 1).

According to this algorithm, the 1-VERSION corpus has been built-up of 338 unarranged text units, the MAIN – of 1 570 text units, patched up into sections of different subject categories (As is seen from A Appendix (Fig. 15)). From the broad list of possible subject categories in the sections, the MAIN corpus focuses on fiction, historical, medical, scientific, sociological literature.

3 The Dictionary of Naturalized Lexical- and Grammatical Information for the Whole List of Unknown Unique Words (File *UNKNOWN.S.dic*)

3.1 “Purity” Check of the Corpus ‘1 million corpus.noc’

To get better results on the task of the Dictionary of Naturalized Lexical- and Grammatical Information creation, it is necessary to realize the more extended “purity” check of the corpus ‘1 million corpus.noc’. From this corpus, with the help of Levenshtein algorithm [2], the search of wordforms with the high (0.8) level of similarity of one wordform to another was realized. As a result, comparing the created dictionaries of known (about 150 000) and unknown (about 50 000) wordforms, the authors have found almost 30 % of similar wordforms, which, as a matter of fact, must belong to known wordforms, despite the fact that the NooJ program has recognized them as unknown.

Below, the general problem points are given in terms of “purity” check of the whole text corpus, which may rather effectively be solved using the abovementioned Levenshtein algorithm:

1. the occurrence of Latin letters in many words in the texts of the Belarusian corpus;
2. dialectal words of the Belarusian language;
3. Russian words;
4. orthographic mistakes;
5. different letter case processing.

Table 1. Descriptive algorithm of the first one-million corpus for the Belarusian NooJ module

№	Action	Result
1	Text material collecting and subject specification for prospective corpora	The draft corpus for the Belarusian NooJ module is made. (The name is 1 million corpus – Very big[noc]). Chosen subjects: without the subject specification. (TEXT TOTAL: 338)
2	<p>The “purity” check of the developed corpus:</p> <ol style="list-style-type: none"> 1. <i>Unknown or unidentified symbols;</i> 2. [In addition]: words with apostrophe (’); 3. [In addition]: words with hyphen/ (–) symbol. <p>2, 3 [All possible word occurrences with <i>apostrophe (’)</i> and <i>hyphen (-)</i> must be searched in texts of the corpus because the NooJ program is unable (in the process of Linguistic Analysis) to parse them correctly (in a context such words are incorrectly lexically divided): e.g. 1) <i>nað’ exaцb</i> (incorrect); 2) <i>цѣнна- ciнi</i> (incorrect). <i>(In both cases the words are divided by the space and NooJ can’t identify them as one unit).</i></p> <p>E.g. 1) <i>nað’ exaцb</i> (correct); 2) <i>цѣнна-ciнi</i> (correct)]</p>	<ol style="list-style-type: none"> 1. There are NO ANY unknown or unidentified symbols in the given corpus (all texts must be encoded to UTF-8). 2. 5818 wordforms with apostrophe (’) are found in the given corpus. (The special NooJ-grammar SearchWordFormsWithApostrophe(FIRST VERSION).nog has been applied in this case) 3. In the given corpus we detected: <ul style="list-style-type: none"> – 25 615 occurrences of all wordforms with hyphen/ (–) symbol (The special NooJ-grammar SearchWordFormsWithHyphen.fst has been applied in this case); – 19 728 occurrences of all unique (1 occurrence per match) wordforms with hyphen/ (–) symbol (The special NooJ-grammar SearchWordFormsWithHyphen.fst has been applied in this case); – 23 663 occurrences of all wordforms with hyphen/ (–) symbol using special <WORDFORMWITHHYPHEN> query. (After the realization of Linguistic Analysis throughout the corpus, applying necessary NooJ resources general_be.nod [4, 5] and SearchWordFormsWithHyphen.nog); – 18 351 occurrences of all unique (1 occurrence per match) wordforms with hyphen/ ‘-’ symbol, using special <WORDFORMWITHHYPHEN> query. (Such results were got after the realization of Linguistic Analysis throughout the corpus, applying necessary NooJ resources general_be.nod and SearchWordFormsWithHyphen.nog).
3	The wordforms total is counted in the developed corpus	Wordforms (ALL) (<WF> = 1 884 971) Wordforms (ALL, unique [1 occurrence per match]) (<WF> = 197 712)
4	The unknown words search is realized from the developed corpus	Unknown words (ALL) (<UNK> = 140 235) Unknown words (ALL, unique [1 occurrence per match]) (<UNK> = 50 186)
5	Creation of the dictionary of unknown (<UNK>) unique word usage [* .dic]	File UNKNOWNS.dic (49 749 unknown words (after the necessary correction of the unique <UNK> = 50 186 result))
6	Creation of the dictionary of naturalized lexical- and grammatical information for the whole list of unknown unique words (on the basis of the file UNKNOWNS.dic)	[The more extended “purity” check of the developed corpus. (The “disposal” of Roman alphabet letters and other problematic cases in texts of the corpus)]

7	The known words search is realized from the developed corpus	Known words (ALL) (<DIC> = 1 750 447) Known words (ALL, unique [1 occurrence per match]) (<DIC> = 148682)
8	Realization of the developed corpus time consuming operation of lexical information selection	The whole dictionary is created on the basis of the given corpus: <i>1 mil c_BE.dic</i> (197 712 unique occurrences)
9	The punctuation marks search is realized from the developed corpus	Punctuation marks (ALL) (<P> = 605 512) Punctuation marks (ALL, unique [1 occurrence per match]) (<P> = 61)
10	The digrams search check is realized from the developed corpus	<u>Digrams</u> <i>Concordance</i> (Digrams = 1 787 305 [occurrences]) <i>Dictionary</i> : (<i>1 mil cor VB_(digrams).dic</i>)
11	The Part of Speech tagging of words is realized within the developed corpus	<u>Parts of Speech</u> (<NOUN> = 523 193) → ALL (<NOUN> = 58 376) → unique [1 occurrence per match] (<VERB> = 326 739) → ALL (<VERB> = 45 682) → unique [1 occurrence per match] (<ADJECTIVE> = 194 682) → ALL (<ADJECTIVE> = 38 218) → unique [1 occurrence per match] (<ADVERB> = 172935) → ALL (<ADVERB> = 3 869) → unique [1 occurrence per match]

These issues are also solved in the NooJ program, though it takes far more time and effort, because the program has to process a large scope of information.

3.2 Statistical Analysis of the Text Corpus ‘*1 million corpus.noc*’

The following steps have been taken at this stage:

1. The corpus ‘*1 million corpus.noc*’ Linguistic Analysis (As is seen from A Appendix (Fig. 16)).
2. The search of wordforms (all wordforms, which are present in the corpus) using special queries (**<WF>**, **<UNK>**, **<DIC>**, **<NOUN>**, **<VERB>**, **<ADJECTIVE>**, **<ADVERB>**).
3. Export of the matches into text files (*.txt).
4. Text files with exported data are stored in a special database, where the unique wordform clustering was realized and the number of wordforms was counted.

3.3 Machine-Learning Algorithms Application for the Part-of-Speech Tagging [3] of Unknown Wordforms

1. The **main attribute of a wordform** for the Part-of-speech tagging process with the machine-learning algorithms, specified for the purpose, was *three ending letters* of each wordform in the dictionary of unknown wordforms. (As is seen from Fig. 1.)

Row ID	S Col0	S Col1	S *Col2	S Col3	I Col4	I Col5
Row33	/Aksionov_SamobeldzajBeliba_AL...	102 — каля 6,5 мілья геста...	асягза...	выядрана 294 тыс. геставай зямлі Беларусь	1918	В
Row34	/Aksionov_SamobeldzajBeliba_AL...	нак. геставай, 2 асягза...	тыт	геставай зямлі, Беларусь — краіна лясоў	1946	В
Row35	/Aksionov_SamobeldzajBeliba_AL...	вазлеў лясу ў гайнак рэ...	Гранца	Дзятр, Сок адносна да зямлі	2044	В
Row36	/Aksionov_SamobeldzajBeliba_AL...	лесе ў пэйзаж рэч Прыпяць...	Дзятр	Сок адносна да зямлі рэчыцкага	2053	В
Row37	/Aksionov_SamobeldzajBeliba_AL...	у гайнак рэч Прыпяць, Дз...	Сок	адносна да зямлі рэчыцкага забаронена	2060	В
Row38	/Aksionov_SamobeldzajBeliba_AL...	рэчыцкіх асяродках, рэ...	асягза...	распространены і пачытаемы мугацкі... 3	2275	В
Row39	/Aksionov_SamobeldzajBeliba_AL...	путэй... 38. Чарнобыль...	Сс	7, 24, 49, 101, 149. У гэтых асяродках, 29 красавіка 19...	2374	В
Row40	/Aksionov_SamobeldzajBeliba_AL...	загісторыі ў Польшчы...	Румыні	30 красавіка — у Швейцарыі і Гішпаніі	2521	В
Row41	/Aksionov_SamobeldzajBeliba_AL...	Германіі, Мусоліні, Румыні...	Швейцары	і Гішпаніі 29 ліпня, 1—2 мая — у	2547	В
Row42	/Aksionov_SamobeldzajBeliba_AL...	Красавіка — у Швейцарыі...	Італі	1—2 мая — у Францыі, Бельгіі, Нідэрландах	2549	В
Row43	/Aksionov_SamobeldzajBeliba_AL...	нае — у Францыі, Бельгіі...	Венецыя	паўночнай Грэцыі, 3 мая — у Італіі	2623	В
Row44	/Aksionov_SamobeldzajBeliba_AL...	Валодарыцкай, найчымэй...	Італі	Італіі, Турцыі... Засячыты на вяселле	2644	В
Row45	/Aksionov_SamobeldzajBeliba_AL...	паўночнай Грэцыі, 3 мая...	Турцы	Турцыі... Засячыты на вяселле	2673	В
Row46	/Aksionov_SamobeldzajBeliba_AL...	Грэцыі, 3 мая — у Італіі...	Турцы	Засячыты на вяселле вышэйшага палкадзейна	2682	В
Row47	/Aksionov_SamobeldzajBeliba_AL...	загісторыі ў Яўцы, 4...	Італі	5-га — у Італіі, 5-6 мая	2824	В
Row48	/Aksionov_SamobeldzajBeliba_AL...	у Італіі, 5-6 мая — у...	Італі	1-га мая. Мая да тыдня	2843	В
Row49	/Aksionov_SamobeldzajBeliba_AL...	16 мая — у Італіі	Італі	Мая да тыдня стартуваў, аб	2847	В
Row50	/Aksionov_SamobeldzajBeliba_AL...	наступнага чэрвеня...	Мінск	Мінска аднавілі выкладкі Спартыўнага вучэлішча па	3006	В
Row51	/Aksionov_SamobeldzajBeliba_AL...	Мінска аднавілі выкладкі Сп...	разлічэн...	1992 г. С. 82. Чварцеры рэагента, нумары	3054	В
Row52	/Aksionov_SamobeldzajBeliba_AL...	выкладкі Спартыўнага вуч...	С. 82.	Чварцеры рэагента, нумары аб	3074	В
Row53	/Aksionov_SamobeldzajBeliba_AL...	Спартыўнага вучэлішча па р...	С	82. Чварцеры рэагента, нумары аб'ектаў	3077	В
Row54	/Aksionov_SamobeldzajBeliba_AL...	разлічэнні, 1992 г. С. ...	навучны	аб'ектаў "Сювітаў", па-ранейшаму	3116	В
Row55	/Aksionov_SamobeldzajBeliba_AL...	С. 82. Чварцеры рэагента...	або	"Сювітаў", па-ранейшаму зберажэн	3117	В
Row56	/Aksionov_SamobeldzajBeliba_AL...	навучны аб'ектаў "Сювітаў",	разлічэн	зберажэн і сёння скарэктаваны	3118	В

Fig. 1. The excerpt of the NLP system database with data from NooJ

- The following algorithms are being applied:
 - Decision Tree;
 - Clustering;
 - Neural Network.
- Firstly, the dictionary [5, 6] of known wordforms was downloaded into the system¹. This dictionary was meant to “train” all possible word paradigms by the above-named algorithms. In other words, the algorithms’ *learning* is realized, and 30 % of known wordforms were taken for *its* realization. (As is seen from Fig. 2.)

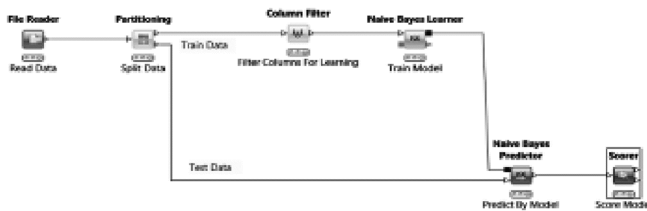


Fig. 2. Train model

Then, 70 % of checked remaining data were realized by the derived model of machine-learning algorithms to verify the proficiency of the model estimated as rather high.

- After that, the existing dictionary of unknown wordforms (UNKNOWN.S.DIC) was “passed” through the machine-learning model. The results produced rather high degree of correct assignment of unknown wordforms to one or another part of speech.

¹ The Part-of-Speech Tagging process can be realized not only in one particular NLP system but also in many others (including integrated interactive systems), where these three algorithms, mentioned above, can be applied.

And that, even at the elementary level (here, the main wordform attribute for the realization of Part-of-Speech Tagging process, i.e. *three ending letters* of each wordform, is meant), confirms the effectiveness of the given machine-learning model. (As is seen from Figs. 3, 4 and 5.)

Row ID	\$ Col0	\$ Col1
Row1750	адвсечко	чко
Row1751	адвсоблввўся	ўся
Row1752	адвсь	сь
Row1753	адвсю	сю
Row1754	адвсья	сья
Row1755	адвсё	сё
Row1756	адвсск	сск
Row1757	адвсчак	чак
Row1758	адвснтрацыя	цыя
Row1759	адвў	ў
Row1760	адвўнату	нату
Row1761	адб	дб
Row1762	адбрэ	рэ
Row1763	адблнц	лнц
Row1764	адбаранў	анў
Row1765	адбкнў	кнў

Fig. 3. Unknown data

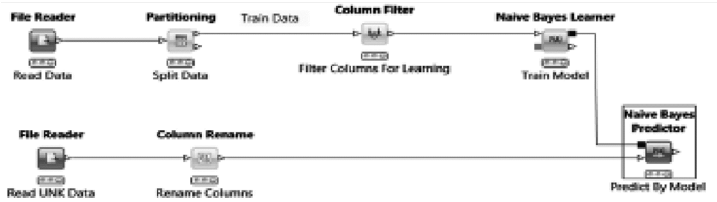


Fig. 4. POS prediction

Row ID	\$ Seq	\$ Ending	\$ PartOfSpeech
Row1750	адвсечко	чко	OTHER
Row1751	адвсоблввўся	ўся	VERB
Row1752	адвсь	сь	OTHER
Row1753	адвсю	сю	OTHER
Row1754	адвсья	сья	VERB
Row1755	адвсё	сё	OTHER
Row1756	адвсск	сск	OTHER
Row1757	адвсчак	чак	OTHER
Row1758	адвснтрацыя	цыя	OTHER
Row1759	адвў	ў	OTHER
Row1760	адвўнату	нату	OTHER
Row1761	адб	дб	OTHER
Row1762	адбрэ	рэ	VERB
Row1763	адблнц	лнц	VERB
Row1764	адбаранў	анў	VERB
Row1765	адбкнў	кнў	VERB

Fig. 5. Predicted POS results

- There are possible variants of data check results by using the aforementioned model (through the example of VERB; other parts of speech are considered as OTHER) (Table 2).

Table 2. Possible variants of data check results

Variants	Actual	Model prediction
True positive	VERB	VERB
True negative	OTHER	OTHER
False positive	OTHER	VERB
False negative	VERB	OTHER

4 Part-of-Speech Tagging Countercheck

The Part-of-Speech Tagging Countercheck on unknown words was realized with the help of Levenshtein algorithm (on basis of the file UNKNOWNNS.dic).

One more task was to work out the dictionary of unknown words usage. The assignment for developers is the maximum reduction of the dictionary sizes and determination of unknown words values for their further correction and inclusion in the dictionary of the one-million Belarusian corpus. (As is seen from Fig. 6.)

№	А	В	С	Д	Е	Ж	З	И	К	Л
10	5	Сяў -- то-е, шло ў гісторыю ўвасходзілі	Сяў -- то-е, шло ў гісторыю ўвасходзілі	Similarity	PartOfSpeech	As is	WV	Grash	Comments	
11	10	ІНШЫЙ	іншы	0,999999999	ADVERB	1	1	1		
11	11	ІНШЫЙ/Ы	іншы/ы	0,999999999	ADJECTIVE	1	1	1		
11	12	ІНШЫЙ/Ы/Ы	іншы/ы/ы	0,999999999	ADJECTIVE	1	1	1		
11	13	ІНШЫЙ/Ы/Ы/Ы	іншы/ы/ы/ы	0,999999999	ADJECTIVE	1	1	1		
11	14	ІНШЫ/Ы/Ы	іншы/ы/ы	0,999999999	NOUN	1	1	1		
11	15	ІНШЫ/Ы/Ы/Ы	іншы/ы/ы/ы	0,999999999	NOUN	1	1	1		
11	16	ІНШЫ/Ы/Ы/Ы/Ы	іншы/ы/ы/ы/ы	0,999999999	NOUN	1	1	1		
11	17	ІНШЫ/Ы/Ы/Ы/Ы/Ы	іншы/ы/ы/ы/ы/ы	0,999999999	NOUN	1	1	1		
11	18	ІНШЫ/Ы/Ы/Ы/Ы/Ы/Ы	іншы/ы/ы/ы/ы/ы/ы	0,999999999	NOUN	1	1	1		

Fig. 6. Unknown words in the summary table of the Part-of-Speech Tagging Countercheck

The main feature of the algorithm applied to the Belarusian one-million corpus is that the algorithm does not change the words in texts after editing, but makes it possible for users to see comments on various mistakes made in the texts incorporated in the Belarusian one-million corpus.

The words included in the dictionary were classified in groups of the unknown for various reasons:

- The words written in the Latin alphabet or having some Latin letters (**a bavyazkov, atrgml_vayutsets, an akhoplepa**);
- The words written by a tarashkevitsa, i.e. substandard spelling which, however, is used by a rather large number of people, especially the Internet users. The existence of alternative spelling is caused by the historical reasons (**абараназдольнасьці, абвешчання, абвясцілі**);
- Words with spelling errors (**абслўгоўванню, абцаюць**);

- Words with recognition errors after scanning (**магілёўскага, встагоддзем**);
- Words of foreign languages (**perfekt, deutsche, eine**);
- Proper nouns (**Дзятлава, Анатоля**), etc.

The main objective at the stage of unknown words recognition is the definition of their morphological characteristics, i.e. assignment of parts of speech value to 49 749 wordforms. The algorithm of Levenshtein revealed parts of speech of unknown words, picked up a possible correct form of the usage, and also gave an index of probability of correct forms. (As is seen from Figs. 7 and 8.)

№	А	В	С	Д	Е	Ж	З	И	К	Л
№	0	Бел -- том. што всправіцца ЛІКОВЫ	Словы -- том. што справіцца ад перапісу	Бел -- том. што справіцца	PartOfSpeech	Адзін	Мног.	Час	Склад	Склад
10	0	магілёўскі	магілёўскі	0,000000000	ADVERB	1	1	1		
11	0	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
12	0	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
13	0	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
14	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
15	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
16	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
17	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
18	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
19	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
20	0	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
21	0	магілёўскага	магілёўскага	0,000000000	FOREIGN	1	1	1		

Fig. 7. Linguists’ checkout of wordforms recognized by Levenshtein algorithm in the summary table of the Part-of-Speech Tagging countercheck

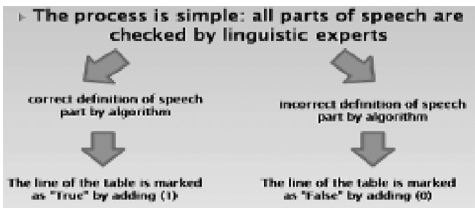


Fig. 8. The process of Parts-of-Speech tagging countercheck

The stage of manual editing is carried out after computer-assisted Part-of-Speech definition, i.e. the algorithm can correctly reveal all parts of speech on formal grounds. The algorithm is simple: all parts of speech are checked by linguistic experts. In case of the correct Part-of-Speech definition by the algorithm this line of the table is marked as “truly” (1). In an opposite case – “lie” (0). (As is seen from Figs. 9 and 10.)

№	А	В	С	Д	Е	Ж	З	И	К	Л
№	0	Бел -- том. што всправіцца ЛІКОВЫ	Словы -- том. што справіцца ад перапісу	Бел -- том. што справіцца	PartOfSpeech	Адзін	Мног.	Час	Склад	Склад
10	0	магілёўскі	магілёўскі	0,000000000	ADVERB	1	1	1		
11	1	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
12	1	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
13	1	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
14	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
15	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
16	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
17	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
18	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
19	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
20	1	магілёўскага	магілёўскага	0,000000000	NOUN	1	1	1		
21	1	магілёўскага	магілёўскага	0,000000000	FOREIGN	1	1	1		
22	1	магілёўскага	магілёўскага	0,000000000	ADJECTIVE	1	1	1		
23	1	магілёўскага	магілёўскага	0,000000000	FOREIGN	1	1	1		

Fig. 9. Editing the results of linguist’s checkout of wordforms recognized by Levenshtein algorithm in the summary table of the Part-of-Speech Tagging countercheck

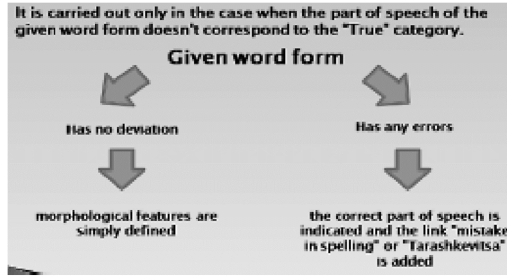


Fig. 10. The process of editing

If the part of speech of a specified wordform doesn't correspond to the validity, the editing stage comes, namely indications of the correct part of speech. If the corresponding wordform has no deviations (without spelling errors, unclear symbols, and also not foreign-language words), in this case morphological features are simply defined. If a word has a wrong spelling, the correct part of speech is indicated and the link “mistake in spelling” is specified. The same happens to the words written in a tarashkevitsa only with another link: “Tarashkevitsa”.

The parts of speech noted by the “NULL” category are mainly proper names and therefore are defined by the algorithm described above: indication of the correct part of speech and assignment to this line of “true” value. (As is seen from Fig. 11.)

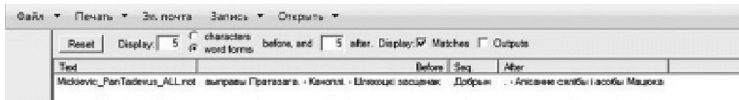


Fig. 11. Addressing to the context

In case the word meaning is not clear or causes doubts, it is necessary to address to a context, namely to the corpus.

At the end of this stage the quantity of unknown words was decreased that allowed to pass to the following stages of the first Belarusian one-million corpus improvements. (As is seen from Fig. 12.)

№	А	Б	В	Г	Д	Е	Ж	З	И
1	а	а	а	а	а	а	а	а	а
2	а	а	а	а	а	а	а	а	а
3	а	а	а	а	а	а	а	а	а
4	а	а	а	а	а	а	а	а	а
5	а	а	а	а	а	а	а	а	а
6	а	а	а	а	а	а	а	а	а
7	а	а	а	а	а	а	а	а	а
8	а	а	а	а	а	а	а	а	а
9	а	а	а	а	а	а	а	а	а
10	а	а	а	а	а	а	а	а	а
11	а	а	а	а	а	а	а	а	а
12	а	а	а	а	а	а	а	а	а
13	а	а	а	а	а	а	а	а	а
14	а	а	а	а	а	а	а	а	а
15	а	а	а	а	а	а	а	а	а
16	а	а	а	а	а	а	а	а	а
17	а	а	а	а	а	а	а	а	а
18	а	а	а	а	а	а	а	а	а
19	а	а	а	а	а	а	а	а	а
20	а	а	а	а	а	а	а	а	а
21	а	а	а	а	а	а	а	а	а
22	а	а	а	а	а	а	а	а	а
23	а	а	а	а	а	а	а	а	а
24	а	а	а	а	а	а	а	а	а
25	а	а	а	а	а	а	а	а	а
26	а	а	а	а	а	а	а	а	а
27	а	а	а	а	а	а	а	а	а
28	а	а	а	а	а	а	а	а	а
29	а	а	а	а	а	а	а	а	а
30	а	а	а	а	а	а	а	а	а
31	а	а	а	а	а	а	а	а	а
32	а	а	а	а	а	а	а	а	а
33	а	а	а	а	а	а	а	а	а
34	а	а	а	а	а	а	а	а	а
35	а	а	а	а	а	а	а	а	а
36	а	а	а	а	а	а	а	а	а
37	а	а	а	а	а	а	а	а	а
38	а	а	а	а	а	а	а	а	а
39	а	а	а	а	а	а	а	а	а
40	а	а	а	а	а	а	а	а	а
41	а	а	а	а	а	а	а	а	а
42	а	а	а	а	а	а	а	а	а
43	а	а	а	а	а	а	а	а	а
44	а	а	а	а	а	а	а	а	а
45	а	а	а	а	а	а	а	а	а
46	а	а	а	а	а	а	а	а	а
47	а	а	а	а	а	а	а	а	а
48	а	а	а	а	а	а	а	а	а
49	а	а	а	а	а	а	а	а	а
50	а	а	а	а	а	а	а	а	а

Fig. 12. The Concatenation-in-Paradigm results

According to the resulting data, the special Concatenation-in-Paradigm list was made after the countercheck of recognized by the Levenshtein algorithm unknown words (previously exported from the NooJ-dictionary file UNKNOWNNS.dic)) in order to create the additional NooJ general_be.nod dictionary. (As is seen from Fig. 13.)

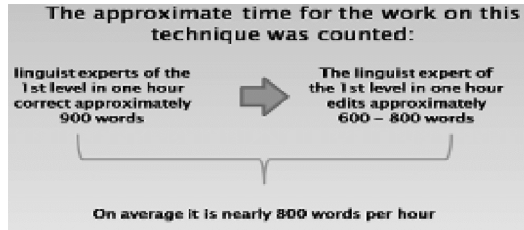


Fig. 13. The approximate time count for the work on the technique

This calculation allows to predict the work on this technique in the future and to estimate degree of overall performance in comparison with other techniques.

5 Comparison of Lexical and Grammatical Base of the Belarusian N-Corpus [6] with Dictionary Properties’ Definition File of the Belarusian NooJ Module

In a similar manner a comparison of lexical and grammatical base of a Belarusian N-corpus with dictionary properties’ definition file of the Belarusian NooJ module was made. The Belarusian N-korpus is the first widely available general Belarusian corpus. The Belarusian N-korpus currently contains ~50,000 texts (~30,000,000 tokens) taken from fiction, newspapers, journals and on-line editions. The texts of the corpus are grammatically annotated and contain metatextual information.

Verb Class (10)						
Класіфікацыя	Belarusian N-Corpus	Category	Belarusian N-adj module	Only in N-corpus	Only in NooJ	In both
Імя	Амалецтва (8)	Number	Singular	0	0	7
	Мноства (9)		Plural			
Рэч	Мрэжчынства (3)	Gender	Мужчынства	2	0	7
	Жанчынства (7)		Жанчынства			
	Імя (5)					
	Мноства (9)					
Час	Пачатак (3)	Tense	Future	1	0	3
	Наступна (1)		Present			
	Канцак (6)		Past			
Загортка	Загортка (8)	Reflexive	Reflexive	1	0	8
	Незагортка (9)					
Ліч	Ліч (1)	Mood	Imperative	0	1	1
	Ліч (1)		Indicative			
Парадкавы	Парадкавы (7)	—	—	3	0	8
	Непарадкавы (1)					
	Парадкавы/Непарадкавы (1)					
Спражэжы	Спражэжы (1)	—	—	1	0	8
	Другае (2)					
Дзеяслова	Дзеяслова (1)	—	—	1	+	1
	Дзеяслова (1)					

Fig. 14. Morphological characteristics of verb

The comparative analysis was performed on the morphological characteristics of different parts of speech listed in dictionaries of both programs. After the structure analysis of both Belarusian N-corpus and the Belarusian NooJ module, it can be concluded that the programs have quite developed system of characteristics of speech parts, but nevertheless some categories need to be improved, what was found out in the process of comparing lexical and grammatical bases. Comparison of the morphological characteristics of Verb is presented on Fig. 14.

6 Conclusion

In conclusion, the first one-million corpus for the Belarusian NooJ module is suitable for research in the following aspects:

1. *Words polysemy processing in texts of different subjects;*
2. *Polysemic punctuation marks processing;*
3. *New lexical items search.*

Besides, the one-million corpus is valuable for solving other important tasks:

- Conduction of several experiments in order to specify the **syntactic and morphological grammar use efficiency** of texts of each subject in the corpus, at minimum as well as maximum level.
- Taking thorough measures in order to create the *subject domain generator*. (This will be then very useful for the formation of special subject-oriented NooJ dictionaries.)
- The usage of the given corpus (in the most extent) in the process of Text-to-Speech synthesis with the help of available programs [7], required for such a process, and also when testing newly created applications.
- Carrying-out of a comparative analysis of this corpus with the same corpora in other languages (taking into account all necessary rules, language features in texts of each current corpus, various possible emerging issues, while building syntactic and morphological grammars, etc.).

Thus, it is essential that the first one-million corpus for The Belarusian NooJ Module has practical application in any line of linguistic research. In the near future the corpus is planned to be expanded up to approximately 5–10 million words.

Acknowledgements. Many thanks to T. Okrut, J. Baradzina, A. Fiodarau for their help in revising the language of this paper.

A Appendix

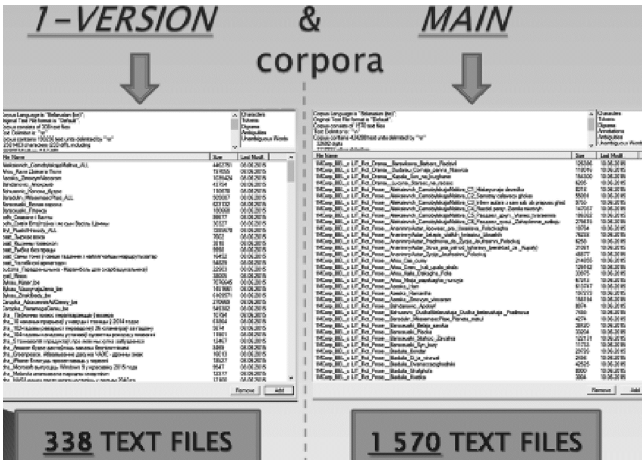


Fig. 15. 1-VERSION and MAIN corpora

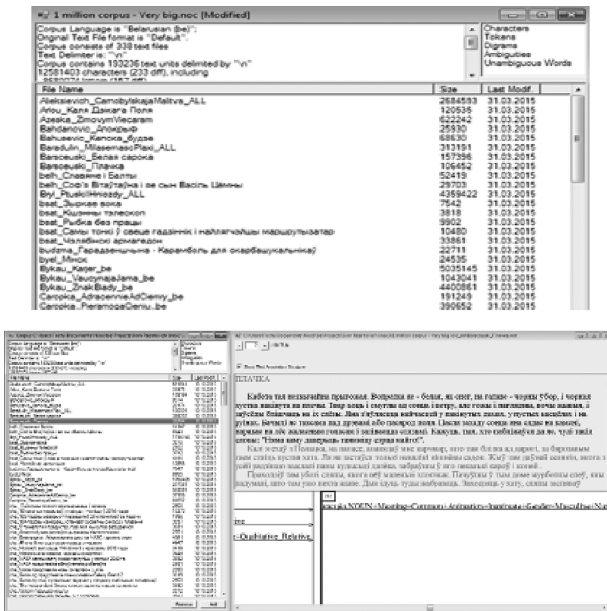


Fig. 16. Realized Linguistic analysis of the 1-VERSION corpus

References

1. NooJ: a linguistic development environment [Electronic resource] (2015). <http://www.NooJ4nlp.net/>. Accessed 08 May 2015
2. The Levenshtein-Algorithm [Electronic resource] (2015). <http://www.levenshtein.net/>. Accessed 24 Sept 2015
3. Taylor, P.: Text-to-Speech synthesis. In: Taylor, P. (ed.) Text Decoding, pp. 89–92. Cambridge University Press, Cambridge (2009). Chapter 5
4. Hetsevich, Y.: Overview of Belarusian and Russian dictionaries and their adaptation for NooJ. Hetsevich, Y., Hetsevich, S. In: Vučković, K., Božo, B., Max, S. (eds.) Automatic Processing of Various Levels of Linguistic Phenomena: Selected Papers from the NooJ 2011 International Conference, pp. 29–40. Cambridge Scholars Publishing, Newcastle (2012)
5. Hetsevich, Y.: Accentual expansion of the Belarusian and Russian NooJ dictionaries. Hetsevich, Y., Hetsevich, S., Lobanov, B., Skopinava, A., Yakubovich, Y. In: Donabédian, A., Khurshudian, V., Max, S. (eds.) Formalising Natural Languages with NooJ : Selected Papers from the NooJ 2012 International Conference, pp. 24–36. Cambridge Scholars Publishing, Newcastle (2013)
6. Аўтаматызаваная апрацоўка сімвальных выказаў у тэкстах для сістэмы сінтэзу беларускага маўлення. Беларускі N-корпус [Electronic resource] (2015). <http://bnkorporus.info/>. Accessed 17 May 2015
7. Corpus.by. Corpus.by [Electronic resource] (2015). <http://www.corpus.by/>. Accessed 08 May 2015