



LingvoDoc: Working with Cognate Analysis

Natalia Kosheliuk

EasyChair preprints are intended for rapid dissemination of research results and are integrated with the rest of EasyChair.

February 11, 2022

LingvoDoc: working with Cognate Analysis¹

Kosheliuk Natalia

ORCID 0000-0002-5833-7971

Ivannikov Institute for System Programming of the RAS, Moscow (Russia)

NKoshelyuk@yandex.ru

Abstract. This article offers an overview of one of the main LingvoDoc programs – the Cognate analysis option intended for modern etymological and dialect distance research. It describes in detail how this option works and how to implement it. The advantages of the Cognate Analysis algorithm are also highlighted.

Keywords. LingvoDoc, data mining, linguistics, cognate analysis

1 INTRODUCTION

Cognate identification is an important thing for identifying the genetic affinity of languages and their dialects. This method allows linguists to draw conclusions about the development of languages over time and obtain new phonetic and etymological data. Today only a small part of languages has been analyzed in terms of their genetic connections. The main reason is that comparative research in linguistics is still based on the personal work of scientists which is very time-consuming. However, in recent years there has been an active development of methods of computer work on the implementation of these tasks.

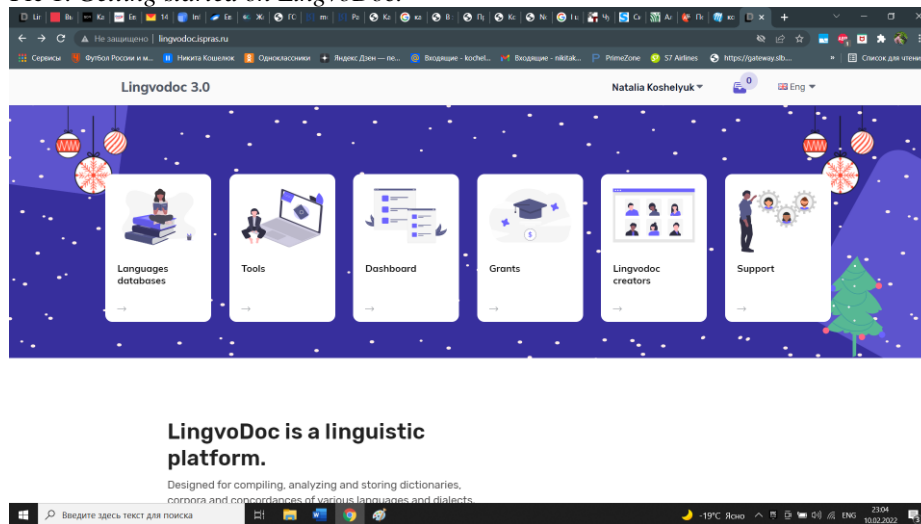
In this paper, using the example of the Mansi corpus, we will show how the Cognate Analysis option is implemented in the LingvoDoc and what its advantage is.

2 HOW TO CONDUCT COGNATE ANALYSIS

Initially, work on the LingvoDoc platform begins with standard authorization, downloading the necessary dictionary, or opening an existing one in the database (Pic. 1).

¹ Supported by Russian Science Foundation, project no. 20- 18-00403 ‘Digital Description of Uralic Languages on the Basis of Big Data’

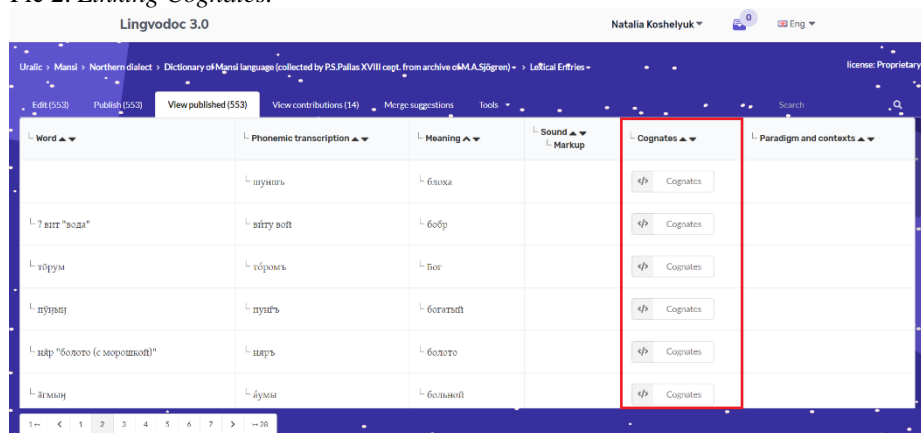
Pic 1. Getting started on LingvoDoc.



To obtain the most accurate and reliable result of etymological analysis, if possible, it is necessary to choose the most qualitative and maximally filled dictionaries (Pic. 2).

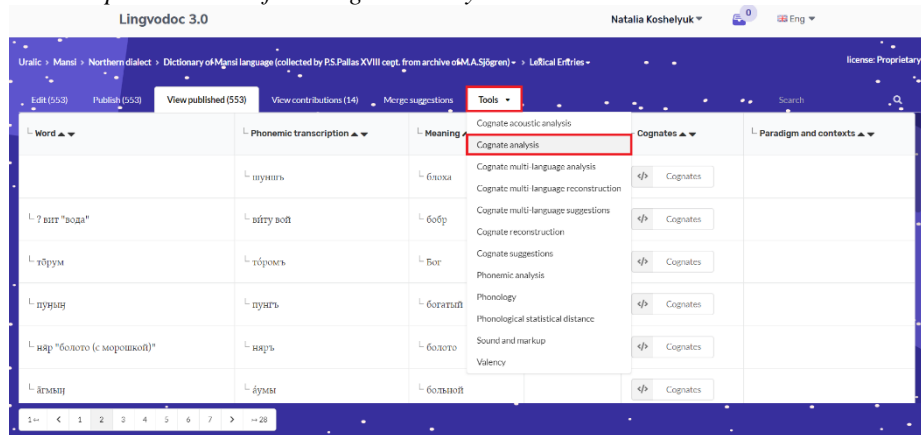
Currently, 23 Mansi dictionaries are connected by etymologies on the Linguistics platform. Most of them are archival texts. The decision to involve them in the analysis is due to the rather high accuracy of the reflection of phonetic oppositions in them. Also, these data are important as one of the first written sources on the Mansi dialects and according to them, you can trace how they have changed over the past 100-250 years.

Pic 2. Linking Cognates.



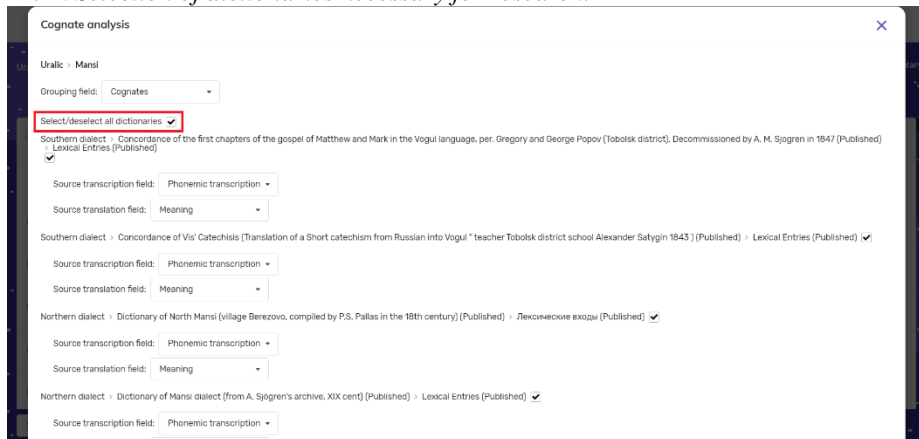
After selecting the necessary linguistic material for Cognate Analysis, you need to go to the "Tools" tab and start the process (Pic. 3).

Pic 3. Implementation of the Cognate Analysis.



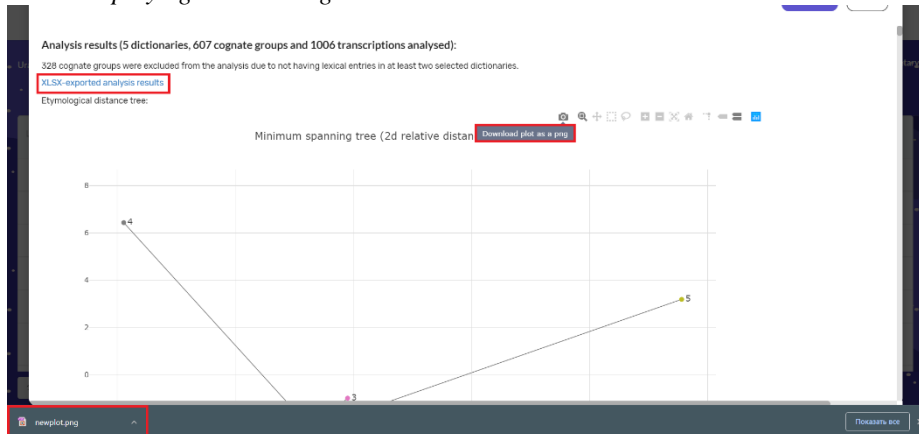
The LingvoDoc platform user also can choose the necessary dictionaries (Pic. 4): you can stop at all available data or limit yourself to at least two dictionaries. Based on our experience, the optimal number of dictionaries subject to processing should not exceed 20. This will be explained by the need to make some editorial changes to combine the rows: manually transferring to reliable – doubtful rows, and to doubtful – single ones. Editorial work with lists of words represented by a large number of columns resulting from the implementation of the Cognate Analysis option is a time-consuming task and can lead to missing important matches.

Pic 4. Selection of dictionaries necessary for research.



After starting the process of Cognate Analysis, which takes from 1 to 5 minutes, the user receives a result that reflects the genealogical distance between the sources involved in the analysis, as well as fully painted rows of reliable and unreliable correspondences of vowels and consonants (Pic. 5-6).

Pic 5. Displaying the resulting result.



Briefly outline the principles of the etymological analysis algorithm on the Linguistics platform: at the first stage, for each character from the transcription proposed in each particular dictionary, the algorithm calculates its correspondences in words from other dialects of the same language associated with etymologies with this dictionary. The main calculation is carried out by the roots of words, based on the fact that in them

the first vowel corresponds to the first vowel, the first consonant corresponds to the first consonant, the second consonant corresponds to the second. Combinations of vowels and vowels, consonants and consonants are possible. At the output, we get a list of correspondences for each pair of idioms.

Pic 6. Displaying the resulting result.

ЭТИМОЛОГИЧЕСКИЙ АНАЛИЗ

1: Dictionary of Mansi language (collected by P.S.Pallas XVIII cent. from archive of H.A.Sjögren) - Lexical Entries (156 форм = 55% от числа соотв.)
 2: The gospel of Matthew and Mark in the Vogul Language (1847-1848) - Lexical Entries (144 форм = 51% от числа соотв.)
 3: Concordance of glossed corpus of Evangel of Matthew 1878 in Konda dialect of Mansi language - Lexical Entries (156 форм = 55% от числа соотв.)
 4: Dictionary of Mansi language, compiled by archpriest Simon Cherkalov (Solkamsk city, 1783) - Lexical Entries (151 форм = 54% от числа соотв.)
 5: Dictionary Upper Pelym dialect of Mansi language, author priest K.Slovcov, 1905 - Lexical Entries (137 форм = 49% от числа соотв.)

Соответствия по начальному гласному:
 Надёжные ряды:
 [a]-[a]-[a]-[a]-[a]
 Сомнительные ряды:
 ? - ? - ? -(i)- ?
 ? -[e]-[e]- ? - ?
 [o]-[o]-[o]-[o]-[a]
 [a]- ? - ? -(e)- ?

Материал - надёжные ряды:

1: Dictionary of M.	2: The gospel of M.	3: Concordance of ...	4: Dictionary of M.	5: Dictionary Upper.
[a]	[a]	[a]	[a]	[a]
аляньть	ачерить	'мороз'	atšerem	'мороз'
аляньть	ач	'а'	av	'а я'
аляньть	ач-очелне	'умирать, погибать'	al-uchv	'погибать-ИВ'
аляньть	ач	atlat	'нет не'	
аляньть	ач	'свиденник'	arkep-	'свиденник'
аляньть	ач-	'собирать'	ach-	'собирать'
аляньть	алянь	'места'	alment-uchv	'места-ИВ'
аляньть	алянь	'ловить'	alisi-achv	'ловить-ИВ'
аляньть	алянь	'другой'	alem	'другой'
аляньть	асрай	'дымвол'	asrai	'дымвол'

The author of the dictionary has the opportunity to download the results of the analysis in Excel format (Pic. 7), analyze, verify the correctness of transcriptions and etymologies that led to non-standard lines of correspondence, and make adjustments to transcription and etymology. Further, the above algorithm can be restarted again already on the material verified by the author.

Pic 9. Lists of correspondences by vowels/consonants, reliable/unreliable lines.

The screenshot shows an Excel spreadsheet with a list of correspondences. The rows are numbered 90 to 116. The first two groups are 'Надежные ряды' (Reliable lines) and 'Сомнительные ряды' (Questionable lines). The first group contains 11 rows (91-102) and the second group contains 15 rows (105-116). The cells contain phonetic symbols and question marks, representing correspondences between different linguistic elements.

Similar to the analysis presented on the LingvoDoc platform, the Excel table displays the names of dictionaries involved in the study, specific examples of words with their translation, as well as identified consonants and vocalisms (Pic. 10-11). According to the fullness of a particular dictionary, it is likely that the file will not be filled with etymologies for all lines.

Pic 10. Results of the analysis: dictionary names and identified etymologies.

11	Материал — надежные ряды:							
12								
13								
14	Dictionary of Mansi language (collect 2; The gospel of Matthew and Mark in 13; Concordance of glossed corpus of Ev 4; Dictionary of Mansi language, compil 5; Dictionary Upper Pelym dialect of Mansi la							
15	[a]	[a]	[a]	[a]	[a]	[a]	[a]	[a]
16	ачерымь	'мороз'	atSerem	'мороз'	Ачерме	'Стужа'	ачеромь	'Мороз, с. 17'
17	амь	'я'	am	'я'			Амь	'Я, с. 27'
18				'я'				
19	влямоть	'убить'	ал-	'умирать, погибать al-uchv	'погубить-INF'			
20			очелне	'погибать, умирать'				
21	птань	'пятьдесят				Ать таль	'Пятилетие'	
22	пть	'пять'						
23	тймь	'нет'		ati	'нет'			
24	втáгмь	'не знаю'		at	'не'			
25			аркоп	'священник'	arker-	'священник'		
26			ак-	'собирать'	ach-	'собирать'		
27			алхлын	'нести'	alment-uchv	'нести-INF'		
28			алислахь	'ловить'	alisl-achv	'ловить-INF'		
29			алым	'другой'	alem	'другой'		
30			асрай	'дьявол'	asrai	'дьявол'		

Pic 11. Results of the analysis: identified phoneme series.

Row	Word	Transcription 1	Transcription 2	Transcription 3	Transcription 4	Transcription 5	Transcription 6	Transcription 7	Transcription 8	Transcription 9	Transcription 10
121	Материал	наблюдение редис:									
122											
123	1. Dictionary of Mansi language (rollet) 2. The gospel of Matthew and Mark in 3. Concordance of glossed corpus of Fv 4. Dictionary of Mansi language. comel 5. Dictionary Upper Pelym dialect of Mansi language, author pri										
124	[m]	[m]	[m]	[m]	[m]	[m]	[m]	[m]	[m]	[m]	[m]
125	мазь	'земля'		та	'страна'	Мя	'земля'	Мя	'земля, с. 12'		
126			машыне	'одежда'		земля'					
127			машыне	'одежда'	та	'одежда'		Мяушь	'Одежда, с. 19'		
128	миньень	'уйди'	менень/мыньень	'идти'	тепу-счv	'идти-INF'					
129					теп-счv						
130	морюх	'морозика'						Морюх пуй	'Морозика ягоды, с. 18'		
131	мáйтимэу	'старуха'					Метъь анкуэу	'Престарелый'			
132	мáйтимэу	'старый'					Мэтъьмэу	'Старуха'			
133	мáйтимáйтэу	'старик'									
134	мáйтэу	'грудь'					Мáйтэу	'Грудь'			
135	мáйтэу	'пенень'					Мáйтэу	'Пенень'			
136	мáйтэу	'ийца'					Мáйтэу	'Йяцо'			
137	мáйтэу	'мáйтэу'			таг	'мáйтэу'					
138			мáйтэу	'немонэу'	таг	'немонэу'					
139			мáйтэу	'одежда'	таг	'одежда'					
140			мáйтэу	'одежда'	таг	'одежда'					
141	[n]	[n]	[n]	[n]	[n]	[n]	[n]	[n]	[n]	[n]	[n]
142	небáсь	бумага	непек	бумага	перек	бумага		Непэк-онса	бумажные деньги, с. 10		
143	конвоймь	пишу'	ханш-	писать'	чанш-	писать'		Наньэ маньагетам	'Грамоте учусь, с. 7		
144								Непэксьуэ	Бумага, с. 4		
145								Наньэ	Газета, с. 7		

Based on the analysis carried out to clarify phonetic and phonological transcription and to calculate regular correlations in words of related languages related by etymologies, at the next stage a program is carried out to determine the proximity of dialects: first, the series of multiple correspondences in the selected group of dictionaries are identified by the vowels of the first syllable and consonants of the first syllable (Pic. 12).

Pic 12. Matrix of proximity of dialects by vowels and consonants of the first syllable.

Row	Header	Col 1	Col 2	Col 3	Col 4	Col 5	Col 6	Col 7	Col 8	Col 9	Col 10	Col 11	Col 12	Col 13	Col 14	Col 15	Col 16	Col 17	Col 18
1800	Соответствия по главному первому слогу (после согласного) (вес: 1):																		
1801																			
1802																			
1803	Dictionary of Mansi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1804	The gospel of Matth	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1805	Concordance of glo	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1806	Dictionary of Mansi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1807	Dictionary Upper Pe	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1808																			
1809																			
1810	Соответствия по начальному согласному (вес: 1):																		
1811																			
1812																			
1813	Dictionary of Mansi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1814	The gospel of Matth	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1815	Concordance of glo	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1816	Dictionary of Mansi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1817	Dictionary Upper Pe	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1818																			
1819																			
1820	Суммарная матрица:																		
1821																			
1822	Dictionary of Mansi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1823	The gospel of Matth	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1824	Concordance of glo	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1825	Dictionary of Mansi	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1826	Dictionary Upper Pe	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1827																			
1828																			

According to Cognate Analysis, based on the identified series of correspondences, all dialects are compared in pairs, the number of transitions

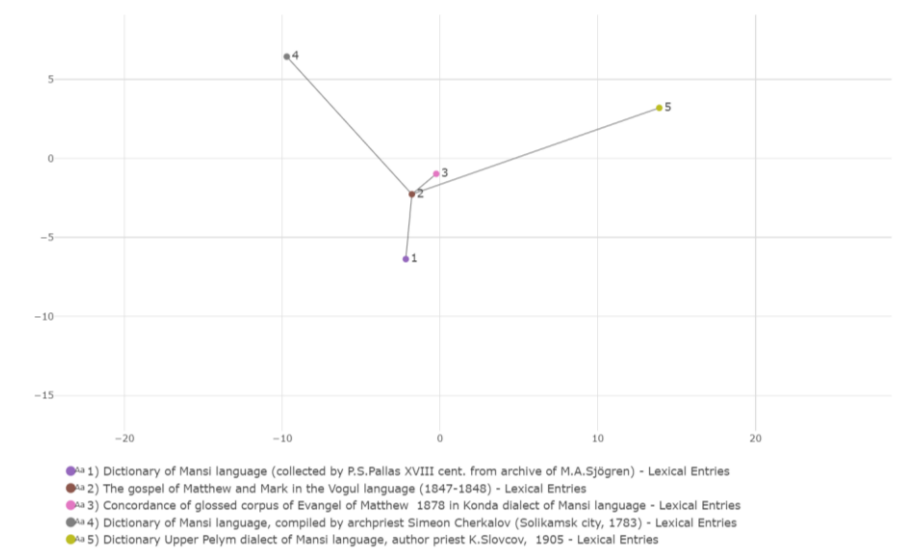
between them is calculated, i.e. a matrix of the analyzed phenomenon is formed: 1) vowels of the first syllable, 2) consonants of the first syllable, 3) sum (Pic. 13).

Pic 13. *The final matrix of dialect proximity.*

1819					
1820	Суммарная матрица:				
1821					
1822	Dictionary of Mansi: The gospel of Matth	Concordance of glo	Dictionary of Mansi	Dictionary Upper Pelym	dialect of Mansi language, author priest K.Slovcov, 1905 - Lexical Entries
1823	Dictionary of Mansi	0	12	26	15
1824	The gospel of Matth	12	0	14	9
1825	Concordance of glo	12	2	0	7
1826	Dictionary of Mansi	26	14	18	0
1827	Dictionary Upper Pelym	15	9	7	19
1828					

The correlation of reliable series of the initial position of a word allows the program to calculate the distance between the dialects represented in dictionaries and distribute them into groups. As a result, an etymological and phonetic 3D-model of the distance between the dialects of the selected dictionaries is constructed. The resulting distance represents the number of similar and different phonetic transitions between languages and dialects (Pic. 14).

Pic 14. *3D-model of the distance between dialects.*



As a result, we get a semblance of a family tree, where the dictionaries involved in the analysis are marked with dots with the appropriate color.

CONCLUSION

As can be seen from the proposed review of the implementation of the Cognate Analysis, this program allows you to process large data arrays in semi-automatic mode to establish etymological and phonological systems of dictionaries presented on LingoDoc. The advantage of working with this option is the ability to quickly and fairly reliably refine and recheck the boundaries of many phonetic phenomena, identify special features and processes that have occurred in languages over a given period of time, and so on.

REFERENCES

1. LingvoDoc Homepage, <http://lingvodoc.ispras.ru/>. Last accessed 02/06/2022