# "LingvoDoc" as a Solution to the Problem of Conservation of Endangered Languages (Using Example of the Mansi Language)

Natalia Kosheliuk

# «LingvoDoc» as a solution to the problem of conservation of endangered languages (using example of the Mansi language)[1]

Kosheliuk Natalia
ORCID 0000-0002-5833-7971
Ivannikov Institute for System Programming of the RAS, Moscow (Russia)
NKoshelyuk@yandex.ru

**Abstract.** This article, using the example of the Russian language, provides an overview of the main problems concerning the endangered languages of Russia, and describes ways to solve them using a linguistic platform LingvoDoc.

**Keywords.** LingvoDoc, data mining, linguistics, endangered languages, conservation, research.
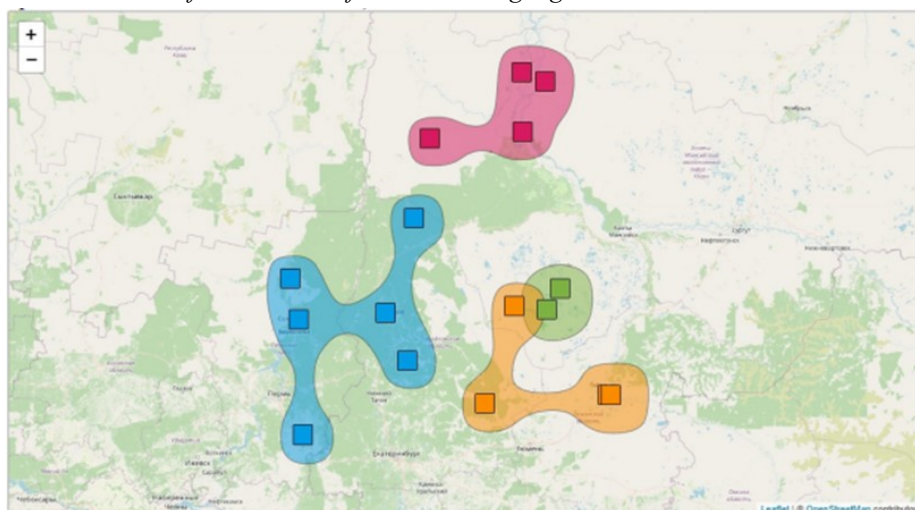
## 1 INTRODUCTION

The Mansi language is the language of native speakers living in the vast territory of Northwestern Siberia. Currently, we are witnessing a catastrophic situation: this language is rapidly dying. According to the annual monitoring of the Institute of Ob-Ugric Research and Development, by 2019, about 1,000 people were still alive who were able to speak the North Mansi dialect, the only one preserved today. These are mostly older people.

It is worth noting that the last speakers of western and southern dialects were recorded in the first decades of the XX century, and the latest data on the eastern Yukonda dialect was recorded in 2013 during an expedition to the Khanty-Mansi Autonomous district by an em-
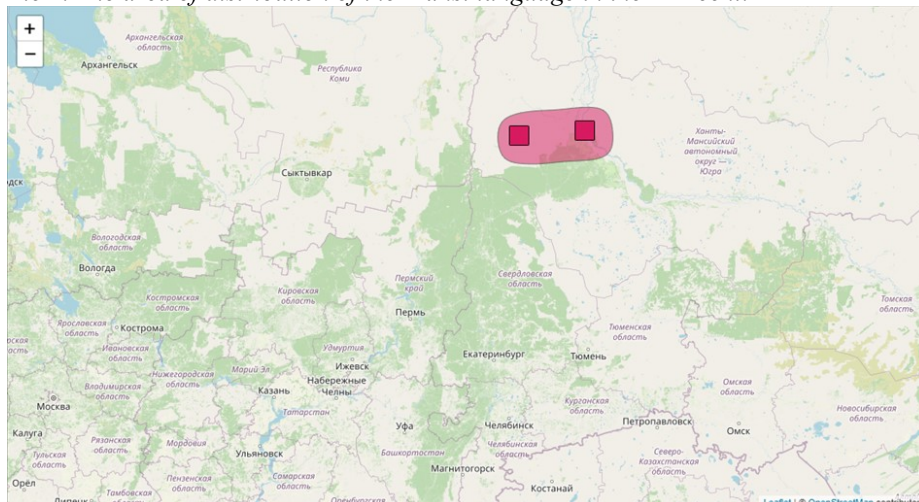
---

ployee of the Institute of Linguistics of the Russian Academy of Sciences. In the village of Shugur, only 2 native speakers were found, Selivanova Elizaveta Ivanovna (89 years old) and Shivtorov Maxim Semenovich (74 years old). Elizaveta Ivanovna passed away 5 days after the departure of the linguist from Shugur, and Maxim Semenovich – in 2018. Thus, the eastern Mansi dialect disappeared literally before our eyes[2] (Pic. 1-2).

*Pic 1. The area of distribution of the Mansi language in the XVIII-XIX cent.*



---

*Pic 2. The area of distribution of the Mansi language in the XXI cent.*



Another problem concerning the Mansi language, as well as other endangered languages, is the difficulty of conducting scientific expeditions. It often turns out that it is often problematic for researchers to get to the places of residence of native speakers: Mansi places of residence are difficult to access for expeditions and require significant resources – physical and financial investments. It is not uncommon for a situation when upon arrival it turns out that the Mansi of those places do not really speak their own language.

This also raises another important problem concerning many languages of Russia – the limited language corpus and, in general, the inaccessibility of most materials for linguists, historians, ethnographers and other specialists. It is known that at present many archival philological documents over the past decades have been tabulated (in particular, the first Mansi Cyrillic translations, which were digitized in the archives of St. Petersburg in 2018-2019. currently, for this reason, they are no longer available) or lost due to transfer to foreign universities (for example, as happened with the archive of A. P. Dulzon). Another point that needs to be mentioned is that many researchers are afraid to publish their materials, transfer them to electronic form, as there are worries about a possible data leak.
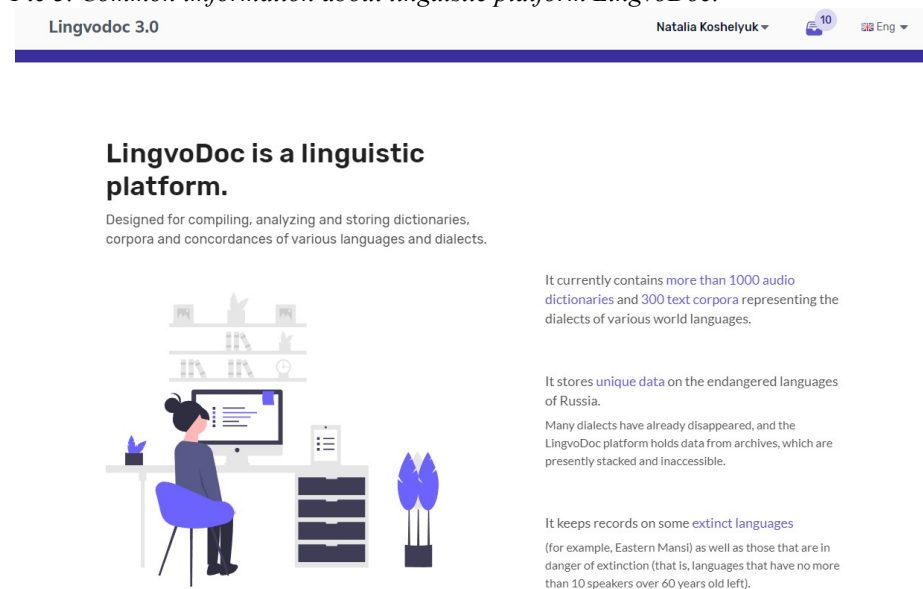
Thus, if no active actions are taken to solve these problems, in the very near future the cultural heritage of Russia in the form of valuable linguistic and archival material may simply disappear.

## 2 HOW LINGVODOC CAN SOLVE THE PROBLEMS OF ENDANGERED LANGUAGES

What solutions does LingvoDoc offer to contribute to the preservation of endangered languages in Russia:

1. The project management (Yu. V. Normanskaya, A. I. Avetisyan, O. D. Borisenko) has a principled position: placing all data in the public domain. Western resources do not practice this – most platforms provide the researcher with data on request.
2. Anyone can register and post their material. Currently, more than 1.000 dictionaries and text corpora have been uploaded to LingvoDoc (Pic. 3).

*Pic 3. Common imformation about linguistic platform LingvoDoc.*
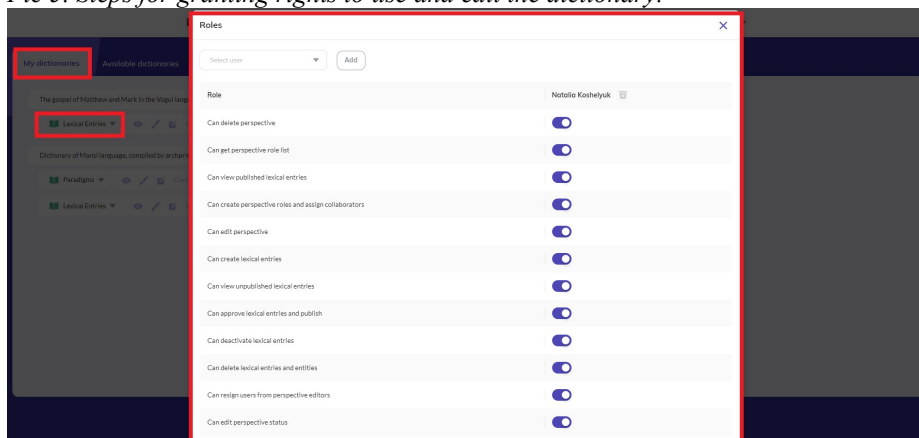
3. If the user believes that the material he is working on at the Lingvo-Doc needs to be improved, it can be hidden in the program settings (My dictionaries -> Hide the dictionary), and published when the author deems it necessary (Pic. 4).

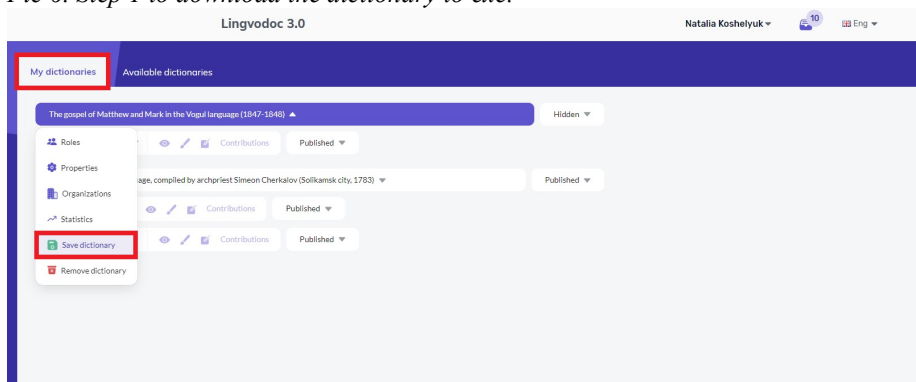*Pic 4. Steps to switch the dictionary to invisible mode.*



4. To prevent plagiarism (Pic. 5), the platform has implemented the option of granting rights to work and edit LingvoDoc resources: you can adjust the list of people who can interact with your corpus or dictionary (My dictionaries -> Lexical entries -> Roles -> Choose the right one).

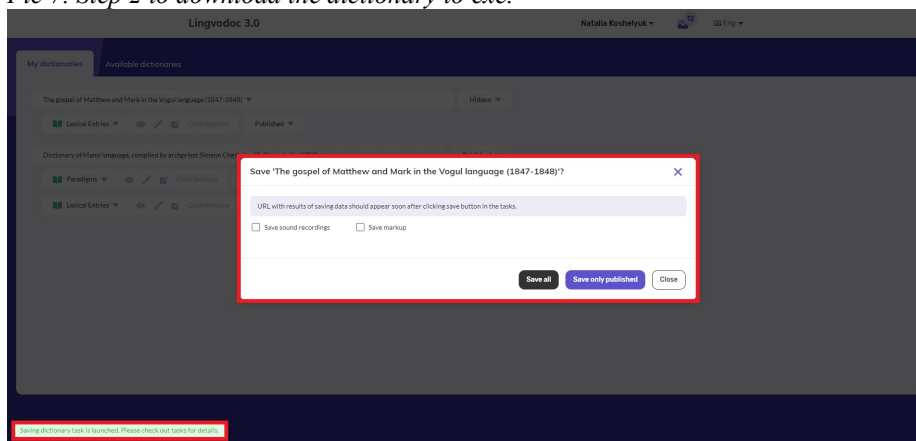*Pic 5. Steps for granting rights to use and edit the dictionary.*

5. Each dictionary or corpus posted on the platform can be registered as IPR on the Gosuslugi website (https://www.gosuslugi.ru – for Russia). This legally secures the copyright to the posted data (Pic. 6-9).
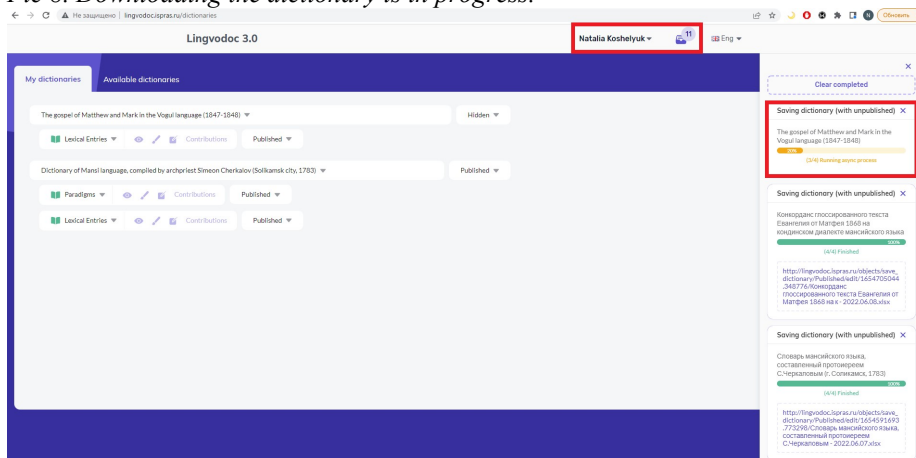
*Pic 6. Step 1 to download the dictionary to exe.*



*Pic 7. Step 2 to download the dictionary to exe.*

*Pic 8. Downloading the dictionary is in progress.*

*Pic 9. Downloaded dictionary in exe.*

It is also necessary to mention another opportunity provided by Linguodoc – the publication of monographs as an indicator of the publication activity of various scientists and scientific institutions.

For today, the amount of data that has already been posted on the LingvoDoc (more than 1,000 sources) allows us to reach the level of publishing monographs from 3-5 books per year for one author. This is ensured by the published amount of data on the LingvoDoc and its processing by a special algorithm that allows you to convert drafts of a future monograph in pdf format: a list of related words, the editing of which can eventually serve as the basis of an etymological dictionary with an average volume of 300-500 pages (Pic. 10).

*Pic 10. A draft of the future etymological Mansi dictionary based on one source posted on the platform.*

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| | A15080 | | | fx | | | | |
| 15057 | | | | | | Словарь чердь | коталъ | 'солнце' |
| 15058 | | | | | | Словарь чердь | уадъ | 'ветер' |
| 15059 | | | | | | Словарь верхо | хоталъ | 'солнце' |
| 15060 | | | | | | Словарь верхо | хоталъ | 'день' |
| 15061 | | | | | | Словарь верхо | вотъ | 'ветер' |
| 15062 | | | | | | Словарь верхо | вотъ | 'вихрь' |
| 15063 | | | | | | Словарь карпи | Уот | 'Ветер' |
| 15064 | | | | | | Словарь карпи | Коатль | 'Солнце' |
| 15065 | | | | | | Словарь карпи | Коатль | 'День' |
| 15066 | | | | | | Словарь карпи | Котлиетъ / Кот | 'Полдень' |
| 15067 | | | | | | Конкорданс М: | котылъ | 'день' |
| 15068 | | | | | | Конкорданс пе | хöтэл пакапа | 'на востоке' |
| 15069 | | | | | | Конкорданс пе | хöтэл пакэпан: | 'на востоке' |
| 15070 | | | | | | Конкорданс пе | хöтэл | 'день' |
| 15071 | kɔ‚åt (P, 387) | палъ катъ палъ | Рука | | | Словарь по па | киул(-) | 'подчинение' |
| 15072 | | | | | | Словарь салы | chütlel | 'mane (утро)' |
| 15073 | | | | | | Словарь сосьв | wo:rtip ka:t | 'левая рука' |
| 15074 | | | | | | Словарь нары | kvark | 'плечо' |
| 15075 | | | | | | Словарь нары | kvɛ par | 'плечо 2' |
| 15076 | | | | | | Словарь нары | kva par | 'плечо, плечи 1 |
| 15077 | | | | | | Материалы из | къпътаръ | 'плечо' |
| 15078 | | | | | | Нарымский ди | кöтпар | 'плечи' |
| 15079 | | | | | | Материалы П.( | когонбаръ | 'плечо' |
| 15080 | | | | | | Материалы из | когондопаръ | 'плечо' |
| 15081 | | | | | | | | |

In the old days, it took linguists decades to carry out such work – collecting material, linking each lexeme with an etymological link and subsequent editorial work on one book. A review of this feature clearly demonstrates how such tasks can now be handled in a few weeks and months.

## CONCLUSION

The opportunities provided by the LingvoDoc for the preservation of endangered languages are an example of how each of the linguists can work with their own languages, as well as quickly publish their articles and monographs. And for higher educational institutions and scientific institutes, this is an example of how each of them can significantly increase the indicators of their institutions in terms of publication activity.

## REFERENCES

1. LingvoDoc Homepage, http://lingvodoc.ispras.ru/. Last accessed: 24/06/2022