# Survey of Natural Language Processing for Vertical Mongolian:
# Current Situation

## Sargai Yunshaab
November 2019

Contents

1.  The state of Natural Language Processing for vertical Mongolian

Mongolian is a language spoken, written and read by about 5.5m people, of whom some 3m live in Mongolia, a democratic, independent state in eastern Inner Asia, while about 2.5m live in the "Inner Mongolian Autonomous Region" (IMAR), a province-level administrative area in the People's Republic of China. Mongolian is agglutinative, with suffix chains in the verbal and nominal domains, and has subject–object–predicate as its basic word order. In Mongolia, the Cyrillic script is used for writing Mongolian, whereas in the IMAR an older script is used, which is written vertically. Also known as the *Hudum Mongol bichig*, this script is derived from the Old Uyghur script and is written in vertical lines top-down, right to left. For foreign-language or mixed-language texts such as this report, the script is often printed horizontally.

Vertical Mongolian natural language processing research started in the 1980's, considerably later than many other languages. In 30 years of development, the basic corpus, grammatical information dictionary, editing and typesetting systems, and common office software have been developed and put into use.[1] Researchers at Inner Mongolia University in Hohhot, the capital of the IMAR, have completed phrase boundary recognition and phrase-structure relationship determination,[2] as well as shallow parsing, such as predicate automatic identification.[3] In terms of code conversion, S. Loglo proposed a code conversion method based on dictionaries and rules, which converts non-national standard code to international standard code, with a conversion accuracy of 98.19%.[4] In terms of spelling proofreading (spellchecks), Hua Shabao has implemented a rule-based spelling error discrimination system, which can distinguish root errors, additional component errors and vowel collocation errors without correction.[5] S. Loglo further proposed a Mongolian proofreading method based on uncertain finite

---

[1] S. Loglo, "Design and realization of a graphical editing system for Mongolian dependency treebank, " *2016 International Conference on Asian Language Processing (IALP)*, Tainan, pp. 176-179, 2016.
[2] Hua Shabao, Dabhurbayar, "A phrase-tagging research in Mongolian corpus", *Journal of the Central University for Nationalities (Philosophy and Social Sciences Edition)*, vol. 33, no. 5, pp. 64-67, 2006.
[3] Serguleng Wang, D. Sarana, Nasunurtu, "Design and realization of automatic annotation for modern Mongolian predicate segment", *Proceedings of 11th national symposium on minority languages*, pp. 420-427, 2007.
[4] S. Loglo, "Research on the General Algorithm of Mongolian Code Conversion [J]." *Journal of Inner Mongolia University (Philosophy and Social Sciences)* 2 (2009).
[5] Hua S. "Modern Mongolian automatic proofreading system–MHAHP. J. Inner Mongolia Univ." *Philos. Soc. Sci. Ed* 4 (1997): 49-53.

automata.[6] The method uses finite state automata to organize Mongolian dictionaries, improves the operation speed of the algorithm and corrects the Mongolian words with correct font but wrong pronunciation. However, manual intervention is needed to select the correct pronunciation. In terms of morphological processing of Mongolian, Nashunuritu combined the characteristics of Mongolian word-endings to form a dictionary and used rules to segment suffixes in Mongolian.[7] Hou Hongxu et al. put forward a Mongolian word segmentation method combining statistics and rules.[8] This method uses a Mongolian statistical language model as the basis for disambiguation of multiple candidate results, effectively improving the accuracy of segmentation. Bao et al. used the method of affix segmentation to improve the recognition rate of large-scale Mongolian speech recognition.[9] Zhao Jiandong et al. added a lookahead mechanism to the historical model, which improved the accuracy of part-of-speech tagging in Mongolian.[10]

The main work carried out on Named Entity Recognition (NER) for Mongolian has been as follows: Tonglaga analyzed personal names in Mongolian and their national, temporal and regional characteristics, along with the changing rules governing the internal model regarding people's names, which he used to build a corpus of personal names.[11] A personal name recognition system based on the maximum entropy model was realized, obtaining an F value=89.61 in 1,040 open name tests. Using the conditional random field model, Cai Jingjing constructed seven sets of different feature templates and ran a test of human name recognition, obtaining an F value of 92.64 in a closed test.[12] Wu Jinxing et al. developed a method for place-name recognition by combining the conditional random field method with use of a dictionary look-up for analyzing the composition of place names in Mongolian.[13]

However, the limited available resources and lack of a public annotated corpus have restricted progress in research and development. In addition, Mongolian's word formation mode is different from those of Chinese and English, and there are some problems, such as inconsistency in coding, which have held back the development of Mongolian NER and limited further research on text understanding, machine translation, public opinion analysis, knowledge map construction, intelligent Q & A, etc.[14]

2.   Problems facing Mongolian NER

In order to investigate the current state of vertical Mongolian NER, interviews were carried out in Inner Mongolia, China between August and October in 2019. The interviewees were from Mongolian information-processing teams at the Inner Mongolia University, the Inner Mongolia Publishing House, the Inner Mongolia Daily, and leading Mongolian information technology companies. I collected literature before and after the interviews to support my findings, and drew, for the majority of this report, on Chinese-language papers which I have summarised and translated into English (these are cited in the footnotes). These sources indicated the following difficulties facing NER for vertical Mongolian:

---

[6] Loglo, S. "A Proofreading Algorithm of Mongolian Text Based on Nondeterministic Finite Automaton." *Journal of Chinese Information Processing* 23.6 (2009): 110-115.

[7] Nashunwuritu, "An Automatic Words Segmentation system of the Mongolian root and stem", *Journal of Inner Mongolia University*, vol. 29, no. 2, pp. 53-67, 1997.

[8] Hou, Hongxu, et al. "Mongolian word segmentation based on statistical language model." *Pattern recognition and artificial intelligence* 22.1 (2009): 108-112.

[9] Bao F., Gao G., Yan X., et al. Segmentation-based Mongolian LVCSR approach[C]. In: Proceedings of Acoustics, Speech and Signal Processing (ICASSP), 2013 IEEE International Conference on. IEEE, 2013: 8136-8139.

[10] ZHAO, Jian-dong, Guang-lai GAO, and Fei-long BAO. "Research on History-based Mongolian Automatic POS Tagging." *Journal of Chinese Information Processing* 27.5 (2013).

[11] 通拉嘎. 基于蒙古文语料库的人名自动识别[D]. 北京: 中央民族大学, 2013.

[12] 才晶晶. 基于 CRF 的蒙古文人名自动识别[D].呼和浩特: 内蒙古大学, 2016.

[13] 吴金星, 丽丽, 杨振新.CRF 和词典相结合的蒙古文地名识别研究[J].计算机工程与科学, 2016, 38(05):1046-1051.

[14] Wang W., "Research on Mongolian name entity recognition" (蒙古文命名实体识别研究), Ph.D. thesis, 2018, Inner Mongolia University (in Chinese).

- There isn't a large enough vertical Mongolian corpus of named entities, coding is not uniform or standardized, and there are many wrong words in the existing corpora, such as those with an incorrect font (glyph) or the wrong pronunciation. Compared with English, Chinese and other developed languages, the information technology foundation for vertical Mongolian is quite basic, and depends upon a named entity corpus that is too small, leading to difficulties for research into NER. For a long time, there have been different coding methods used in Mongolian, and this has made it more difficult to collect and organize a corpus. The typography used in the corpus also brings difficulties for compiling statistics and analysing Mongolian corpora[15].

- Vertical Mongolian has no capital letters, unlike the Latin alphabet. This makes it more difficult to judge the boundary of Mongolian named entities, unlike many languages where the initial capital is a clear boundary marker for a named entity. For example, the English or German name "Einstein" is written in traditional Mongolian as "ᠠᠢᠨᠰᠲᠠᠢᠨ", but this does not have an initial capital, making it difficult to judge its boundary.[16]

- The Mongolian lexicon is large because of the characteristics of Mongolian word formation. Unlike English and Chinese, Mongolian words are completed by affixing suffixes, and this process generates different words from the same stem. This method of word formation leads to a large vocabulary in Mongolian, which increases the number of unknown words in the corpus and makes it difficult to train and test the model.[17]

- The grammar order in Mongolian, subject+object+predicate, makes it difficult to recognise the boundaries of a named entity, as a named entity can act as subject or as an object in any sentence.

- The forms in which named entities appear in of Mongolian are changeable and varied. For example, there are clear differences between Mongolian names and Chinese names because of their different characteristics. Chinese names written in Mongolian usually have the same number of words as the number of syllables in the equivalent Chinese names. For example, the Chinese name 李纪恒 (Li Jiheng) is "ᠯᠢ ᠵᠢ ᠾᠧᠩ", which is composed of three words, whereas the Mongolian personal name agula "ᠠᠭᠤᠯᠠ" is one word. Another difference is that some Mongolian names use the patronym or father's name/surname, such as "ᠵ · ᠪᠠᠲᠤᠮᠦᠩᠬᠡ".[18]

- In Mongolian, the phenomenon of compound words is very common. As in other languages, some common nouns, adjectives and even verbs in Mongolian may act as personal names, and the occurrence of compound words is more frequent. For example, the adjective "ᠴᠡᠴᠡᠨ" (meaning "smart") is a common Mongolian name. It is impossible to determine its category by depending only on the word itself, so it must be judged according to context.[19]

- Another problem is that new entity names keep popping up. Like other languages, Mongolian is a dynamic language. It creates new words constantly and absorbs words from other languages to realize its own development. These new words bring more named entities into the lexicon. For example, "ᠾᠸᠯᠢᠩᠭᠧᠷ ᠱᠢᠨ᠎ᠠ ᠲᠣᠭᠣᠷᠢᠭ" (Helingeer New District) is a new place name, which cannot be found in the current corpus. This can bring difficulties for classification.[20]

---

[15] Wang W., "Research on Mongolian name entity recognition" (蒙古文命名实体识别研究), Ph.D. thesis,2018,Inner Mongolia University (in Chinese).

[16] *Ibid.*

[17] *Ibid.*

[18] *Ibid.*

[19] *Ibid*.

[20] *Ibid*.

3.  Currently available products developed by the Inner Mongolia University Computer Science Lab

The Computer Science Lab at the Inner Mongolia University in Hohhot offers a number of directly related products and tools. This information below reproduces the claims made by the Computer Science Lab on its websites about its products (see Figures 1 and 2).



*Figure 1: Product page of the Computer Science Lab at Inner Mongolia University (October 2019)*
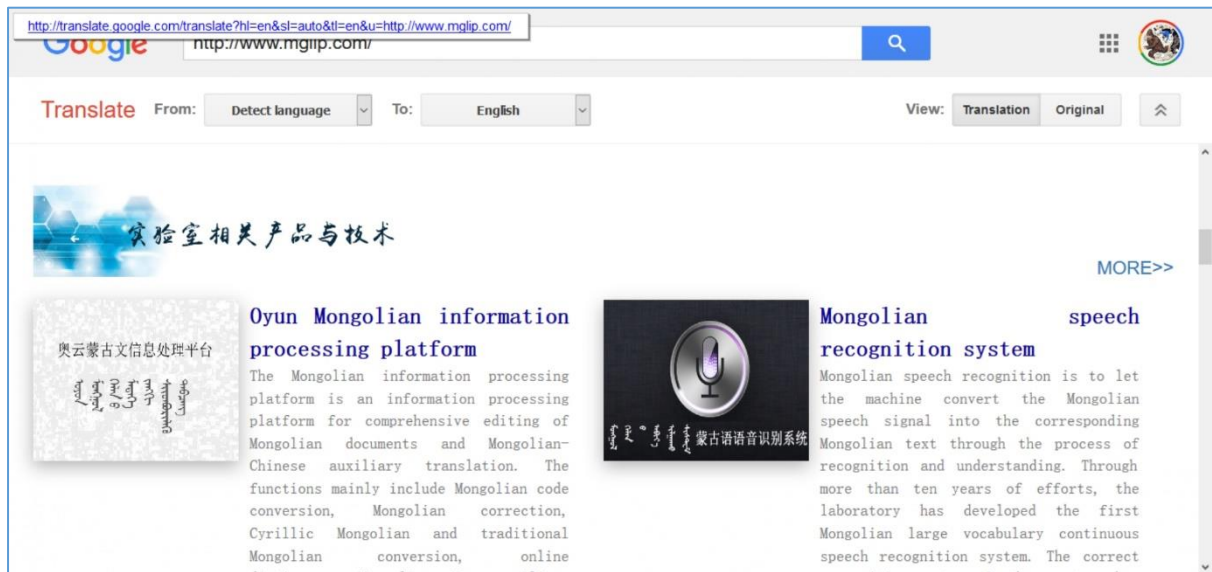
*Figure 2: Google Translation of two items from the Computer Science Lab's product page*

Almost all of these tools are inaccessible to outsiders – they can only be accessed by members of the university – and the site, apart from the home page, appears to be no longer accessible, so I was unable to assess the claims by the Science Lab about most of their products

- The Mongolian Information Processing Platform. This is an information processing platform for comprehensive editing and for auxiliary translation of documents from Mongolian to Chinese. Its functions mainly include Mongolian code conversion, Mongolian correction, conversion between Cyrillic Mongolian and traditional (vertical) Mongolian, an online dictionary, Mongolian-Chinese auxiliary translation, and so on. According to the site, this can meet most of the document information processing needs in Inner Mongolia.
- The Mongolian Print Recognition (OCR) System. [21] According to the site, Mongolian print recognition is a process of transforming Mongolian characters in graphic formats into editable text format by using pattern recognition, artificial intelligence, digital image processing and other technologies. The site claims that recognition accuracy of the Mongolian OCR system is high, and it provides users with a simple and fast editing function in terms of recognition results. This helps users easily solve the digital problems of Mongolian books, newspapers and pictures.
- The Mongolian Speech Synthesis System. The main problem of Mongolian speech synthesis is how to transform visual Mongolian text information into spoken Mongolian, i.e., converting it into data or information that can be heard by the human ear. The Mongolian Speech Synthesis System adopts industry-leading technology, with what is described as a high level of fluency, naturalness and intelligibility.
- The Aoyun Correction System for Mongolian. In order to solve the problem of statistics and retrieval of words with the same appearance but wrong coding, and to speed up the process of developing Mongolian information technology, the laboratory developed a Mongolian proofreading system based on a combination of dictionary look-up, rules and statistics. The system can better correct the words with the same visual form but the wrong coding, and correctly select words with the same visual form but different coding, based on analysis of context.
- The Mongolian Speech Recognition System. Speech recognition in Mongolian aims to let machines transform Mongolian speech signals into the corresponding Mongolian text

---

[21] H. Zhang, H. Wei, F. Bao and G. Gao, "Segmentation-Free Printed Traditional Mongolian OCR Using Sequence to Sequence with Attention Model," *2017 14th IAPR International Conference on Document Analysis and Recognition (ICDAR)*, Kyoto, 2017, pp. 585-590.

through the process of recognition and understanding. As a result of over ten years' efforts, the laboratory has developed the first large-vocabulary continuous speech recognition system for Mongolian in China. The correct recognition rate of the system is said to be more than 90%.

- The Aoyun Mongolian Network Information Retrieval and Monitoring Platform. The platform is a system for news "hot spot" detection, tracking, and content retrieval for Cyrillic Mongolian websites and traditional Mongolian websites. Reportedly, it can automatically grab content, carry out keyword detection on Mongolian websites, and provide a series of functions such as information retrieval, warehousing, statistics, etc., which provides what is described as a convenient service for Mongolian news detection and tracking.

- The System for conversion between Cyrillic Mongolian and traditional Mongolian. The laboratory adopts the method of combining rules and statistics to research and develop a mutual conversion system between Cyrillic Mongolian and traditional Mongolian, with a claimed accuracy of more than 90%. It is claimed to be of great significance for the development of Mongolian science, culture and education.

- The Mongolian Code Conversion System. The laboratory has developed a Mongolian code conversion system using advanced technology. The system provides a mutual conversion capability between Mongolian code and national standard code, which is simple in operation and fast in conversion. It has further accelerated the popularization and use of Mongolian standard coding. This is said to play a very important role in the further application and promotion of historical materials.

- The Mongolian and Chinese machine translation system. The laboratory has established a large number of Mongolian-Chinese dictionaries and corpora. It has also established a machine translation system for Cyrillic Mongolian and traditional Mongolian into Chinese using a combination of rules and statistics. Through this system, content in Mongolian can be filtered and analysed at different levels to help users retrieve Mongolian documents with specific information.

4. Findings from testing the Cyrillic to Vertical Mongolian converter

I was able to access and obtain one tool, the System for conversion between Cyrillic Mongolian and traditional Mongolian. As noted above, vertical (traditional) Mongolian and Cyrillic Mongolian are Mongolian scripts used respectively in China and Mongolia. However, a large portion of Cyrillic Mongolian words can be written in more than one way in vertical Mongolian, which makes conversion from Cyrillic Mongolian to vertical Mongolian difficult. To overcome this difficulty, Bao F et al.'s paper proposed a language model-based approach, which takes advantage of contextual information. Experimental results show that, for Cyrillic Mongolian words that have multiple correspondences in traditional Mongolian, the accuracy rate for this approach reaches 87.66%, thereby greatly improving overall system performance.[22] In this pilot project, our team tested the online converter extensively, as shown in Figure 3.

We found the following problems in the conversion results:

- When Cyrillic is written in all capital letters, the converter doesn't recognize the words.
- The converter doesn't recognize text abbreviations well.
- The converter sometimes does not recognize English letters that are used in Cyrillic Mongolian (e.g. km, kg).
- The converter makes many mistakes when it comes to prepositions (connecting words).

[22] Bao F., Gao G., Yan X., Wang H. (2013) "Language Model for Cyrillic Mongolian to Traditional Mongolian Conversion." In: Zhou G., Li J., Zhao D., Feng Y. (eds), *Natural Language Processing and Chinese Computing. Communications in Computer and Information Science*, vol 400. Springer, Berlin, Heidelberg.

- The converter makes mistakes when there are foreign names since the foreign names don't exist in Mongolian vocabulary lists.
- The converter sometimes transliterates Cyrillic words inaccurately, mistaking it for a word with a similar pronunciation but a different meaning (and sometimes such words are obsolete or do not make sense).
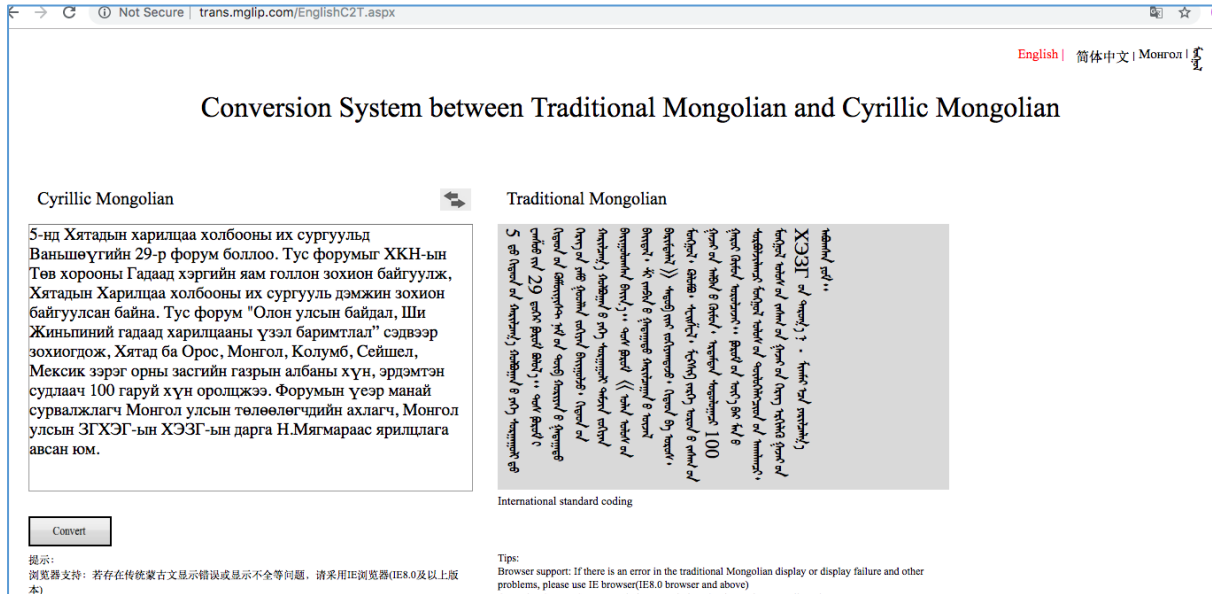


*Figure 3: An example of a test result of conversion from Cyrillic to vertical Mongolian script*

5. Findings from testing the Mongolian Print Recognition (OCR) system

I was also able to run a test of the Mongolian Print Recognition (OCR) System. The tests showed some difficulties for the system in recognizing images of vertical Mongolian script:

- The OCR system only recognizes certain types of fonts (Mongolian Black, Mongolian White, Mongolian Title, Mongolian HaWang, Mongolian Sonin).
- The degree of tilt of the content in the uploaded image affects the accuracy of recognition.
- The background of the image affects the accuracy of recognition.
- The OCR system doesn't recognize hand-written documents (see Figure ), including calligraphed texts (Figure ).
- The OCR makes mistakes on the same word many times even when the source image is very clear (Figure 6, Figure ).
- The OCR system makes mistakes when the contents are written in different fonts in the same image, even for the Mongolian White font and Mongolian Title font (Figure ).

Overall, from the tests I ran, the accuracy of the OCR system appears to be relatively low, at around 40% or lower.

6. Conclusion

NLP is well developed for Cyrillic Mongolian script, but Vertical Mongolian script presents many difficulties for NLP tool developers. So far these difficulties cannot be fully resolved by using a Cyrillic-Vertical Script converter, as the conversion systems still need further development. A large number of

tools for vertical Mongolian script have been developed at Inner Mongolia University, but for the most part these appear not to be accessible to those outside the university or outside the country. More research is needed to establish whether any of these tools could be made available for public use.



*Figure 4: An example of an OCR test result of a hand-written document*



*Figure 5: An example of an OCR test result of a calligraphy document*



*Figure 6: An example of an OCR test result of an image containing clear lyrics*

*Figure 7: An example of an OCR test result of an image containing a slogan*



*Figure 8: An example of an OCR test result of an image containing words written in different fonts*

## Appendix

<u>Major Translated Websites in Vertical Mongolian</u>

List of major websites in vertical Mongolian script translated or adapted from parallel Chinese-language websites. These sites could be used for developing a parallel Chinese-vertical Mongolian corpus

| Website name (Chinese) | Websites in vertical Mongolian (URL) | Source website in Chinese |
|---|---|---|
| 人民网 | http://mongol.people.com.cn | http://www.people.com.cn |
| 中国蒙古语广播网 | http://www.mongolcnr.cn | http://www.cnr.cn |
| 新华网 | http://www.nmg.xinhuanet.com/mg/ | http://www.xinhuanet.com |
| 中国共产党新闻网 | http://mongol.people.com.cn/306956/index.html | http://cpc.people.com.cn |
| 兴安日报 | http://xam.mgyxw.net | http://www.xingandaily.cn |
| 内蒙古科技信息网 | http://mengwen.nmnet.com.cn | http://www.nmsti.com/2014/ |

| 中国蒙古学信息网 | http://mgl.surag.net | http://www.surag.net |
|---|---|---|
| 内蒙古大学蒙古研究中心 | http://mgxzx.imu.edu.cn/mgwb/m1.htm | http://mgxzx.imu.edu.cn |
| 乌兰察布蒙古语新闻 | http://uq.mgyxw.net | http://www.wlcbnews.com |
| 赤峰市政府网 | http://mgl.chifeng.gov.cn/U_index.html | http://m.chifeng.gov.cn |
| 内蒙古自治区政府网 | http://mgl.nmg.gov.cn/index.html | http://www.nmg.gov.cn/ |
| 新疆日报蒙文版 | http://mongol.xjdaily.com | http://www.xjdaily.com.cn |
| 中国蒙古语新闻网 | http://www.mgyxw.net/ | |
| 新疆蒙古语广播网 | http://www.xjmglr.com | |
| 通辽日报 | http://tl.mgyxw.net | |