QR – Issue          QR – Article

**Rustam Abdurasulovich Karimov**
Bukhara State University
A senior teacher of
Foreign Languages Faculty

# PROBLEMS OF LINGUISTIC TAGGING OF UZBEK-ENGLISH PARALLEL CORPUS TEXTS

*Abstract: the article analyzes topic such as: the parallel corpus usually consists of texts on a specific topic. Parallel corpus of state and international official, diplomatic content make up the majority. The most widely used of these is the Corpus of European Documents. The strictness of the translation of such texts is strictly required, and the quality of translation is high. Accuracy in translation makes it easier to equate the two language corpus units. Therefore, in order to align the text of the parallel corpus to give a good result, it is necessary to pay special attention to the quality of the translation when selecting the text for the parallel corpus.*

*Key words: parallel corpus, diplomatic content, quality of translation, language corpus units, linguistic sources and methods.*

*Language: English*

*Citation: Karimov, R. A. (2020). Problems of linguistic tagging of Uzbek-English parallel corpus texts. ISJ Theoretical & Applied Science, 06 (86), 137-140.*

*Soi: http://s-o-i.org/1.1/TAS-06-86-26     Doi: crossref https://dx.doi.org/10.15863/TAS.2020.06.86.26*
*Scopus ASCC: 1203.*

### Introduction

During the pilot phase, the Russian-Chinese parallel corpus received 4 "Reports on the work of the Government of the People's Republic of China" for 2012-2015 and their translations into Russian. At the next stage, it is planned to introduce 6 more official and diplomatic documents. This means that in the first stage it will be possible to parallel the corpus with less material and in the next stage to multiply them. Today, the volume of the corpus consists of 931 paragraphs, 116,668 texts, of which 46,190 text forms are in the Russian language section, and 70,478 in the Chinese language section.

### II.Literature review

Y.M.Mukhin, I.Yan consider the following factors when choosing a source:

1. To improve the quality of the translation, the translated text must be translated by a well-known, qualified professional or official organization. The official website of the Government of the People's Republic of China[1] served as a source of material for this corpus. The materials of this site are translated by the Central Compilation & Translation Bureau; which indicates that the quality of the translation is guaranteed.

2. The selected material must be structurally and semantically stable. For example, the Prime Minister's annual "Speeches" have a repetitive word / sentence / combination, or even text-structure. This is natural. Such duplicate units are important for the translation of the original and the translated text as well as for the automatic marking at a later stage.

3. In contrast to the translation of a work of art, the translation of official-administrative, diplomatic, administrative documents must preserve the content and structure of the original text. Therefore the sequence of such original and translated text sentences / phrases will be the same; structurally mutually identical; will be similar. In this way, the task of aligning text units in two languages is simplified.

### III.Analysis

---

[1] http://cn.theorychina.org/

Principles of corpus marking. The experience of several parallel corpuses, including the Russian-Chinese parallel corpus, was studied in formulating the principles of corpus marking.

In order to get an idea of the stage of parallel corpus correction, we first give an example of parallel texts marked from Russian-Chinese parallel corpus:

1) original text (a):

a1[在财政收支矛盾较大的情况下，我们竭诚尽力，]@‖@  a2[始终把改善民生作为工作的出发点和落脚点，] @| a3[注重制度建设，] @‖@ a4 [兜住民生底线，] @‖@ a5[推动社会事业发展。]

2) translation text (b):

б1[При наличии довольно крупных противоречий между финансовыми доходами и расходами мы со всей искренностью] ‖@ б2[неизменно брали за исходную точку и конечную цель всей своей работы улучшение народной жизни,] @| б3[уделяя особое внимание институциональному строительству,] @‖@ б4[не допуская выхода за нижний предел обеспечения народной жизни] @‖@ б5[и стимулируя развитие социальных сфер.][2].

In this example, elementary discursive units (EDBs) are separated by square brackets. The letter and number between them represent Chinese clauses as well as the syntactic unit in Russian and their order. The vertical line in front of the claus (the "|" sign) indicates the amount in the tree structure to which the claus belongs; The "@" sign indicates the central position of the elementary discursive unit relative to the clauses. As can be seen from the fragment, the Sino-Russian parallel corpus is an elementary discursive unit; types of relationships between them; discursive binding; their semantic characteristics, the central role of the EDB; including other information about its hierarchical structure.

As can be seen from the fragment, the Chinese-Russian parallel corpus is an elementary discursive unit; types of relationships between them; discursive binding; their semantic characteristics, the central role of EDB; including other information about its hierarchical structure.

Below we take a closer look at the layout parameters.

1. *Logical-semantic relationship in the discursive corpus.* Unlike a syntactic relationship, a discursive relationship is primarily logical-semantic

in nature. In his doctoral research, Lee cites a classification that includes 4 groups of relationships (including 17 types), taking into account the importance of asymmetric relationships and discursive connections in the classification. This classification is based on Chinese syntactic theories. In his article, Mukhin cites the types of discursive relationships identified by Lee[3]:

1) parallel relationship: connected, serial, progressive, alternative, comparative;

2) contradictory attitude: contradictory and proportional;

3) causal relationship: private causal, purposeful, conditional, hypothetical, concluding;

4) expanding attitude: descriptive, summarizing, illustrative and evaluative[4].

This relationship classification has been used not only in the Chinese text layout but also in the Chinese-English parallel corpus layout[5].

It is this classification that was used to implement the platform layout created by V.Feng[6]. Similarly, in particular in the formation of the Chinese-Russian parallel corpus, the authors relied on this discursive relationship.

In such a corpus, the marking is done manually on the basis of a special program. The layout platform interface was developed by Chinese scientist V.Feng. On the basis of this interface it is possible to fill the corpus with parallel texts, make the necessary markings, align the units. The multilingual platform can encode different fonts without any hassle. The first area of the platform is dedicated to the task of marking, statistics, marking unification (comparison of markings performed by different users). The second field forms a directory for downloading text and storing marked data. The third area is for data generation, storage, deletion, and editing. At the top of the field is the markup relationship, while in the other field is the original and translated text; a separate field lists the discursive link, the type of discursive relationship, and other parameters.

Marked Russian-Chinese texts are stored in XML files. Such a platform can perform both marking and equalization of a discursive structure at the same time: splitting two parallel texts into a discursive unit, equalizing a discursive structure, defining a discursive relationship, and so on. The presence of such information collected / stored in XML format in the corpus will be useful when performing a search in the corpus later.

[2] «Доклад о работе правительства КНР», 2014 г.
[3] М.Ю.Мухин, Ян И.  Проект создания китайско-русского параллельногокорпуса официально-деловых текстов с дискурсивно-структурной разметкой / Bulletin of the South Ural State University. Ser. Linguistics. 2016, vol. 13, no. 4, pp. 23-31.
[4] Li, Y. Building a Chinese Discourse Corpus with Connective-driven Dependency Tree Structure / Y. Li, W. Feng, J. Sun et al. // In Proceedings of the 2014 conference on Emporical Methods

in Natural Language Processing. – Doha: Qatar, 2014. – P. 2105–2114.    –    http://emnlp2014.org/papers/pdf/ EMNLP2014224.pdf
[5] Feng, W. Alignment and Annotation of Chinese-English Discourse Structure Parallel Corpus / W. Feng // Journal of Chinese Information Processing.– 2013. –27 (6).– P.158–164. – http:// jcip.cipsc.org.cn/CN/abstract/abstract1795.shtml
[6] Ўша манба.

Observing the experience of creating a parallel corpus, we came to the conclusion that the study of the criteria for selecting text on the corpus, the principles of their marking and alignment, as well as the specifics of the software, this process should be gradual. Based on the existing software base, you can easily create a similar parallel corpus.

The presence of a translation dictionary in a parallel corpus further expands its possibilities and serves as a convenient lexicographical tool.

### IV.Discussion

N. Gochev's views on the bilingual / translated dictionary, which will be part of the parallel corpus, are noteworthy[7]. In the article "Corpus of parallel Russian and Bulgarian texts" N. Gochev emphasizes the importance of the parallel corpus and describes its components as follows[8]: "The Russian-Bulgarian parallel corpus consists of two parts. The first section (the corpus itself) is a database at the level of the text adapted / aligned at the sentence level, containing Russian and Bulgarian parallel texts. The second section (Russian-Bulgarian parallel corpus translation dictionary) - the lexical unit being translated is sorted in the corpus from the typological point of view; in which, as a dictionary article, the equivalent corresponding to each unit for translation and the full explanation given to it; the rule of composing a text in the target language based on the nature of the target language; the influence of the original language on the translated language, the transformation of translation; there is information in the form of an exhortation article about errors and inaccuracies in the translation.

In another article, G.N. Gochev discusses the second part of the Russian-Bulgarian parallel corpus - the article structure of the translator-dictionary[9].

In the Russian-Bulgarian parallel corpus, windows such as "Lexical Equivalent", "Translation Transformation", "Incorrect Translation and Errors" are the main part of the corpus translator. In this article, the author discusses the general rules for the structure of a dictionary article in a parallel corpus translator-dictionary on the example of the verb *говорить* in Russian.

Before compiling the dictionary article on the verb *говорить*, 3393 infinitive and verb forms of the *говорить* verb encountered in the Russian-Bulgarian parallel corpus were analyzed. Analyzes showed that

of these, 2782 cases were translated with lexical equivalent.

It is known that in the process of translation there are cases of lack of lexical equivalent or inefficient use of the translator (although it is rare). The equivalent has the appearance that it approaches the meaning of the word in Russian, and that they are a unit that does not correspond to each other; such an equivalent is called an invariant equivalent. The results of the analysis show that there are three invariant-equivalents in the verb *говорить*: these are *говоря, казвам, обаждам се.* All of these units are close to the content and ottenka of the Russian verb *говорить*. But it is not possible to use these words interchangeably in any text in translation.

Other equivalents of the verb *говорить* differ significantly in the translated text. They are option-equivalent and are used only when it is possible to exchange them with invariant-equivalent in the text. The main equivalents of the verb *говорить* are *говоря, казвам;* their frequency in the corpus is the basis for such a conclusion: in the corpus the word *казвам* occurs 1382 times, the word *говоря* 945 times, and the word *обаждам се* 9 times. There are 17 variants of the *Казвам* invariant-equivalent that occur 292 times, and 9 variants of the *говоря* invariant-equivalent that occur 154 times; the variant of the word forms *обаждам се* has not been determined.

According to the author, grouping by lexical equivalent is the most convenient, optimal principle in compiling a corpus translator. This approach results in a perfectly developed system of equivalents and a convenient way to show it to the user. This is practical for the translator.

Making a dictionary article for equivalents Put an index in the form of a hyperlink to that word in the annotated dictionary of each language; can be created by attaching a comment to the word via a hyperlink. From 2 to 4 comments from a translation dictionary to a single word in the original can be attached via a hyperlink, i.e., multiple equivalents to a single word may be appropriate, or vice versa.

### V.Conclusion

Each instance in the corpus checks to see how well these equivalents fit; they clearly indicate the extent to which the word is used in the language;

[7] Гочев Г.Н. Русско-болгарский параллельный текстов как источник двуязычной лексикографии *II* II CONGRESO INTERNACIONAL "La Lengua y Literatura Rusas en el espacio educativo international: estado actual y perspectivas" En conmemoraciyn de los 55 acos de ensecanza de la lengua rusa en Espaca, Granada, 8-10 de septiembre de 2010, tomo II, Granada, 2010.

[8] Гочев Н. Корпус параллельных русских и болгарских текстов // Горизонты прикладной лингвистики и лингвистических технологий. Международная научная конференция. − Киев, 2009. − С. 36-37.

[9] Гочев Г.Н. Словарьные эквиваленты в русско-болгарском словаре переводчика (На материале русского глагола *говорить* и его переводов в корпусе параллельных русских и болгарских текстов // Теоритические и методические проблемы русского языка как иностранного в традиционной и корпусной лингвистике. Доклады и сообщение десятого международного симпозиума. МАПРЯЛ − Болгария. Велико-Тырново, 2010. – С. 682-693.

depending on the frequency, the researcher selects one of the equivalents.

In a translated dictionary, a dictionary article is constructed on the basis of grouping the equivalent, taking into account the language structure. This approach is consistent with the purpose of this type of dictionary (to give more, more complete information about the translated word).

**References:**

1. Muhin, M.Jy., & Jan, I. (2016). Proekt sozdanija kitajsko-russkogo parallel`nogokorpusa oficial`no-delovyh tekstov s diskursivno-strukturnoj razmetkoj / *Bulletin of the South Ural State University*. Ser. Linguistics., vol. 13, no. 4, pp. 23-31.

2. Li, Y., et al. (2014). *Building a Chinese Discourse Corpus with Connective-driven Dependency Tree Structure.* In Proceedings of the 20 14 conference on Emporical Methods in Natural Language Processing. - Doha: Qatar, . - P. 21 05-21 14. - Retrieved from http://emnlp20 14.org/papers/pdf/ EMNLP201 4224.pdf.

3. Feng, W. (2013). Alignment and Annotation of Chinese-English Discourse Structure Parallel Corpus / *Journal of Chinese Information Processing*, 27 (6), pp.158-164.

4. (n.d.). Retrieved from http://jcip.cipsc.org.cn/CN/abstract/abstract179 5.shtml.

5. Gochev G.N. (2010). *Pyccko-bolgarskij parallelnyj tekstov kak istochnik dvujazychnoj leksikografii* II II CONGRESO INTERNACIONAL "La Lengua y Literatura Rusas en el espacio educativo international: estado actual y perspectivas" En conmemoraciyn de los 55 acos de ensecanza de la lengua rusa en Espaca, Granada, 8-10 de septiembre de 2010, tomo II, Granada.

6. Gochev, N. (2009). *Korpus parallel`nyh russkih i bolgarskih tekstov*. Gorizonty prikladnoj lingvistiki i lingvisticheskih tehnologij.

Mezhdunarodnaja nauchnaja konferencija. (pp.36-37). Kiev.

7. Gochev, G.N. (2010). *Slovar`nye jekvivalenty v russko-bolgarskom slovare perevodchika* (Na materiale russkogo glagola govorit` i ego perevodov v korpuse parallel`nyh russkih i bolgarskih tekstov. Teoriticheskie i metodicheskie problemy russkogo jazyka kak inostrannogo v tradicionnoj i korpusnoj lingvistike. Doklady i soobshhenie desjatogo mezhdunarodnogo simpoziuma. MAPRJaL − Bolgarija. Veliko-Tyrnovo, pp. 682-693.

8. Marcu, D., & Wong, W. (2002). *A phrase-based, joint probability model for statistical machine translation*. Proceedings of the Conference on Empirical Methods in Natural Language Processing. (pp. 87-99). Philadelphia.

9. Och, F.J., & Ney, H. (2003). *Discriminative training and maximum entropy models for statistical machine translation.* ACL Anthology: [site]. URL: Retrieved from http://acl.ldc.upenn.edu/P/P02/P02-1038.pdf

10. Toutanova, K., Ilhan, H.T., & Manning, C.D. (2003). *Extensions to HMM-based statistical word alignment models* // Proceedings of Empirical Methods in Natural Langauge Processing., (pp. 87-94). Philadelphia.

11. Brown, P.F., Della Pietra, V.J., Della Pietra, S.A., & Mercer, R.L. (1993). The mathematics of statistical machine translation: parameter estimation. *Computational Linguistics*, Vol. 19(2).