



UTILIZING PARALLEL CORPORA FOR ELECTRONIC DICTIONARY DEVELOPMENT: AN EXAMINATION OF THE UZBEK-RUSSIAN LANGUAGES

Sukhrob Sobirovich Avezov

Bukhara State University Lecturer of the Department of

Russian Literary Studies

1990senigama@gmail.com

Abstract

This research study delves into the utilization of parallel corpora, with a particular focus on Uzbek-Russian instances, for the development of electronic dictionaries. It scrutinizes the enhancement of quality and precision in dictionary entries, brought forth by parallel corpora, which offer an expanded and comprehensive context based on the genuine use of language. The research incorporates historical context pertaining to the evolution of electronic dictionaries, coupled with an exhaustive analysis of the seminal contributions from leading scholars in this domain. Furthermore, the paper provides an in-depth understanding of the technicalities involved in employing parallel corpus data in software engineering to create proficient and operative electronic dictionaries.

Keywords: parallel corpora, electronic dictionaries, natural language processing, lexicography, computational linguistics, machine translation, bilingualism, cross-linguistic analysis.

Introduction

Parallel corpora comprise collections of texts wherein each text in one language aligns with an accurate translation in a different language. [1] Within the sphere of lexicography, the Uzbek-Russian parallel corpus signifies an immense potential for the fabrication of electronic dictionaries, positioning it as a precious asset for linguistic research, machine learning, and translation.

The assembly of the Uzbek-Russian parallel corpus necessitates a substantial volume of bilingual texts. Such texts could range from books and articles to subtitles and web pages, amongst other sources that encapsulate identical content in both languages. These texts are accumulated, normalized, and subsequently coordinated either at the level of a sentence or a paragraph. [2]



Main Part

An electronic dictionary refers to a lexicographic resource accessible and operational in digital format. [3] These dictionaries extend across multiple platforms, encompassing mobile applications, websites, desktop software, and embedded tools within word processing systems.

The distinguishing features and advantages of electronic dictionaries encompass:

1. **Portability and Availability:** The usability of electronic dictionaries spans various devices, permitting anytime, anywhere access, thereby enhancing their convenience.
2. **Searchability and Navigation:** Electronic dictionaries typically facilitate swift and user-friendly keyword searches, with some dictionaries even proffering audio and symbol search capabilities.
3. **Upgradeability:** Given their digital nature, electronic dictionaries can be readily and promptly updated to incorporate new lexicon and phrases.
4. **Multimedia Integration:** Electronic dictionaries might incorporate audio-visual files for pronunciation guidance, illustrations, and supplementary usage examples.
5. **Educational Opportunities:** A considerable number of electronic dictionaries introduce learning-enhancing features such as quizzes, games, and tests to bolster vocabulary acquisition.

The following categories delineate electronic dictionaries:

1. **Monolingual Dictionaries:** These refer to dictionaries that provide explanations of words within the same language.
2. **Bilingual or Multilingual Dictionaries:** These dictionaries facilitate the translation of words from one language to another.
3. **Specialized Dictionaries:** This category includes dictionaries that focus on particular subjects or fields, such as medical, legal, technical dictionaries, and so on.
4. **Phraseological Dictionaries:** These consist of common expressions or phrases in a specific language.

The preliminary notions of electronic dictionaries originated in the mid-20th century, concomitant with the progression of computer technology. Nonetheless, the significant leap occurred in the 1970s and 1980s, when experimental endeavors were initiated to convert traditional dictionaries into their electronic counterparts.

One of the pioneering and most celebrated projects was executed by the Linguistic Research Center at Lancaster University, UK. Guided by Jeffrey Leah, the researchers fabricated the inaugural electronic dictionary for the English language, laying the foundation for subsequent ventures in this domain.



The advent of personal computers and the internet in the 90s and 2000s catalyzed the surging popularity of electronic dictionaries. They were facile to disseminate and update, simple to operate, and boasted substantially larger storage capacity compared to their printed counterparts.

Frank Richtener from the University of Munich, a key researcher in this realm, has made significant contributions to the COSMAS project. His work, which encompassed the development of a comprehensive electronic dictionary and text corpus, played a pivotal role in shaping the field of electronic lexicography.

Electronic dictionaries have now seamlessly integrated into our daily routine. They are accessible not only on computers but also on mobile devices, offering numerous supplementary features such as audio pronunciation, example sentences, and interactive language learning exercises.

Professor Jan Müller from the University of Leiden has carried out substantial work in the field of electronic lexicography, including the development of one of the earliest electronic dictionary mobile applications.

Through the utilization of parallel corpora, it is feasible to construct electronic dictionaries grounded in the actual usage of the language, as opposed to a theoretical or synthetic context. Initially, it facilitates the identification and categorization of dictionary entries based on their real-life usage context. Secondly, it offers a more profound and extensive context for each word as they are displayed within the framework of bilingual sentences.

Parallel corpora serve as a potent instrument for fabricating electronic dictionaries that facilitate a more effective and precise learning process of foreign languages. [4] A key advantage in utilizing parallel corpora in the construction of electronic dictionaries lies in the potential to formulate dictionary entries anchored in real-world language usage. It implies that each word or phrase is classified and expounded upon based on its practical usage in texts, rather than its theoretical or abstract meaning. Such an approach considers the context, idiomatic expressions, and cultural nuances that may elude traditional dictionaries.

Moreover, parallel corpora enable each word to be showcased in the context of bilingual sentences, affording a richer and wider-ranging context due to the multitudinous semantic nuances. Users can observe how a word or phrase is applied in diverse situations and styles, thus aiding in a superior understanding of its meaning and application. [5]

Beyond generating more precise and beneficial lexicons for language learners, parallel corpora can also be employed to enhance the quality of machine translation.



[6] Machine translation algorithms can be educated on parallel corpora to effectuate more accurate text translations, taking into account the real-world usage and contextual subtleties of words and phrases.

Operations involving parallel corpora primarily boil down to manipulating data in text format. For instance, Python and the pandas library can be utilized to process and analyze such data.

Assume we possess a table with parallel corpora of the Uzbek and Russian languages. The structure of such a table might resemble the following:

Uzbek	Russian
Salom, dunyo!	Привет, мир!
Yaxshi, rahmat.	Хорошо, спасибо.
...	...

Presuming this table is stored in a CSV file named 'parallel_corpus.csv', we can import this file into a pandas DataFrame and manipulate it using Python.

Here's an illustrative code snippet that imports data and subsequently executes a rudimentary electronic dictionary lookup:

```
import pandas as pd

# Загружаем данные
df = pd.read_csv('parallel_corpus.csv')

# Функция для поиска перевода
def find_translation(word, df):
    result = df[df['Узбекский'] == word]['Русский']
    if len(result) > 0:
        return result.values[0]
    else:
        return "Перевод не найден"

# Ищем перевод слова "Salom"
print(find_translation('Salom', df))
```

In this exemplification, the 'find_translation' function accepts a word and a DataFrame as parameters and yields the translation of the word, if discovered in the Uzbek column.



Please acknowledge that this example is simplified and omits several critical facets of handling text data, such as text preprocessing (e.g., transforming to lowercase, removing punctuation) and addressing ambiguities in translation. These issues can be tackled by employing supplemental natural language processing (NLP) tools and methodologies.

Overall, the adoption of parallel corpora in the development of electronic dictionaries signifies a considerable advancement in the field of lexicography and language learning. They supply more accurate and beneficial resources for language learning, while also establishing a foundation for more advanced natural language processing and machine translation technologies.

Parallel corpora also hold immense value for training machine learning models, including those for machine translation and NLP. An abundance of bilingual data enables the training of models that can comprehend and replicate intricate grammatical structures and idiomatic expressions in both languages.

Conclusion

The Uzbek-Russian Parallel Corpus unfolds enormous potential for the fabrication of more effective and precise electronic dictionaries, and moreover, for advancements in the fields of machine translation and NLP. It symbolizes a novel approach to lexicography, underscoring the significance of genuine language utilization and proffering fresh avenues for language research and the evolution of natural language processing technologies.

REFERENCES

1. Захаров В. П., Богданова С. Ю. Корпусная лингвистика. – 2011.
2. Добровольский Д. О. Корпус параллельных текстов и сопоставительная лексикология //Труды института русского языка им. ВВ Виноградова. – 2015. – Т. 6. – С. 413-449.
3. Mengliev B. R., Hamidovna N. L. Problems of language, culture and spirituality in general explanatory dictionaries of Uzbek language //International Journal of Psychosocial Rehabilitation. – 2020. – Т. 24. – №. 3. – С. 378-385.
4. Nigmatova L., Avezov S. ПРИМЕНЕНИЕ МЕТОДОВ NLP В КОРПУСНЫХ ИССЛЕДОВАНИЯХ: ОСОБЕННОСТИ И ОГРАНИЧЕНИЯ //«УЗБЕКСКИЕ НАЦИОНАЛЬНЫЕ ОБРАЗОВАТЕЛЬНЫЕ ЗДАНИЯ ТЕОРЕТИЧЕСКОЕ И ПРАКТИЧЕСКОЕ СОЗДАНИЕ ВОПРОСЫ" Международная научно-практическая конференция. – 2023. – Т. 2. – №. 2.



5. Avezov S. S. MACHINE TRANSLATION TO ALIGN PARALLEL TEXTS //International Scientific and Current Research Conferences. – 2022. – С. 64-66.
6. Sharipov S. ЛЕКСИКОГРАФИЯ (ТАРЖИМА ЛЕКСИКОГРАФИЯСИ) РИВОЖЛАНИШИНИНГ АСОСИЙ ЖИҲАТЛАРИ //ЦЕНТР НАУЧНЫХ ПУБЛИКАЦИЙ (buxdu. uz). – 2022. – Т. 15. – №. 15.