

ИНФОРМАТИКА COMPUTER SCIENCE

IRSTI 20.23.21

10.51889/2959-5894.2024.87.3.009

A.K. Aitim^{1*}, R.Zh. Satybaldiyeva²

¹International Information Technology University, Almaty, Kazakhstan

²Satpayev University, Almaty, Kazakhstan

*e-mail: a.aitim@iitu.edu.kz

A SYSTEMATIC REVIEW OF EXISTING TOOLS TO AUTOMATED PROCESSING SYSTEMS FOR KAZAKH LANGUAGE

Abstract

The development of automated systems for the Kazakh language has gained significant momentum in recent years, driven by the growing need for natural language processing (NLP) tools tailored to underrepresented languages. This systematic review aims to critically evaluate the existing observational tools and methodologies used in building and improving automated systems for learning the Kazakh language. Through a comprehensive analysis of scientific literature, technical reports, and practical implementations, this review identifies key trends, challenges, and advances in the field. The review highlights various linguistic complexities unique to the Kazakh language, such as its agglutinative nature, vowel harmony, and rich morphological structure, which pose unique challenges to developers. Additionally, the study examines the effectiveness of modern tools, including tokenization, part-of-speech tagging, parsing, and machine translation, in processing Kazakh text. The findings show that despite significant progress, there are still significant gaps in the availability and accuracy of these tools, especially when compared to those available for more widely spoken languages. The review concludes with recommendations for future research and development, highlighting the need for more robust datasets, improved algorithms, and collaborative efforts to further advance Kazakh language data science.

Keywords: kazakh language processing, natural language processing, machine translation, transformer models, Kazakh text classification, computational linguistics, Kazakh language.

Ә.Қ. Әйтiм¹, Р.Ж. Сатыбалдиева²

¹Халықаралық Ақпараттық Технологиялар Университетi, Алматы қ., Қазақстан

²Қазақ Ұлттық Техникалық Зерттеу Университетi, Алматы қ., Қазақстан

ҚАЗАҚ ТІЛІНЕ АРНАЛҒАН АВТОМАТТАНДЫРЫЛҒАН ӨНДЕУ ЖҮЙЕЛЕРІНЕ ҚОЛДАНЫЛҒАН ҚҰРАЛДАРЫНА ЖҮЙЕЛІ ШОЛУ

Аңдатпа

Қазақ тілін автоматтандырылған өңдеу жүйелерінің дамуы соңғы жылдарда айтарлықтай қарқын алды, бұл табиғи тілдерді өңдеу құралдарының жеткіліксіз ұсынылған тілдерге бейімделген қажеттілігінің артуына байланысты. Бұл жүйелі шолу қазақ тілінің автоматтандырылған жүйелерін құру мен жетілдіруде қолданылып жүрген бақылау құралдары мен әдістемелерін сыни тұрғыдан бағалауға бағытталған. Академиялық әдебиеттерді, техникалық есептерді және практикалық енгізулерді жан-жақты талдау арқылы бұл шолу осы саладағы негізгі тенденцияларды, қиындықтарды және жетістіктерді анықтайды. Шолуда қазақ тіліне ғана тән әртүрлі тілдік күрделілік, мысалы, оның агглютинативті табиғаты, дауысты дыбыстардың үндестігі, бай морфологиялық құрылымы, әзірлеушілерге ерекше қиындықтар туғызатыны көрсетілген. Сонымен қатар, зерттеу қазақша мәтінді өңдеудегі токенизацияны, сөз бөлігін тегтеуді, синтаксистік талдауды және машиналық аударманы қоса алғанда, қазіргі құралдардың тиімділігін зерттейді. Нәтижелер айтарлықтай прогреске қол жеткізілгенімен, бұл құралдардың қолжетімділігі мен дәлдігінде, әсіресе кеңірек сөйлейтін тілдер үшін

қол жетімді құралдармен салыстырғанда, әлі де айтарлықтай олқылықтар бар екенін көрсетеді. Шолу қазақ тілін өңдеу саласын одан әрі ілгерілету үшін анағұрлым сенімді деректер жинақтарының, жетілдірілген алгоритмдердің және бірлескен күш-жігердің қажеттілігін атап көрсете отырып, болашақ зерттеулер мен әзірлемелер бойынша ұсыныстармен аяқталады.

Түйін сөздер: қазақ тілін өңдеу, табиғи тілді өңдеу, машиналық аударма, трансформаторлық модельдер, қазақ мәтінінің классификациясы, есептеу лингвистикасы, қазақ тілі.

А.К. Айтим¹, Р.Ж. Сатыбалдиева²

¹Международный Университет Информационных Технологий, г. Алматы, Казахстан

²Казахский национальный исследовательский технический университет имени К. И. Сатпаева,
г. Алматы, Казахстан

СИСТЕМАТИЧЕСКИЙ ОБЗОР СУЩЕСТВУЮЩИХ ИНСТРУМЕНТОВ ДЛЯ АВТОМАТИЗИРОВАННЫХ СИСТЕМ ОБРАБОТКИ КАЗАХСКОГО ЯЗЫКА

Аннотация

Разработка автоматизированных систем обработки для казахского языка в последние годы получила значительный импульс, что обусловлено растущей потребностью в инструментах обработки естественного языка, адаптированных для недостаточно представленных языков. Целью этого систематического обзора является критическая оценка существующих наблюдательных инструментов и методологий, используемых при создании и совершенствовании автоматизированных систем для казахского языка. С помощью всестороннего анализа академической литературы, технических отчетов и практических реализаций этот обзор определяет ключевые тенденции, проблемы и достижения в этой области. Обзор подчеркивает различные лингвистические сложности, уникальные для казахского языка, такие как его агглютинативная природа, гармония гласных и богатая морфологическая структура, которые представляют уникальные проблемы для разработчиков. Кроме того, исследование изучает эффективность текущих инструментов, включая токенизацию, разметку частей речи, синтаксический анализ и машинный перевод, при обработке казахского текста. Результаты показывают, что, несмотря на значительный прогресс, все еще существуют значительные пробелы в доступности и точности этих инструментов, особенно по сравнению с теми, которые доступны для более широко распространенных языков. Обзор завершается рекомендациями для будущих исследований и разработок, подчеркивая необходимость в более надежных наборах данных, улучшенных алгоритмах и совместных усилиях для дальнейшего продвижения области обработки казахского языка.

Ключевые слова: обработка казахского языка, обработка естественного языка, машинный перевод, модели трансформаторов, классификация казахских текстов, компьютерная лингвистика, казахский язык.

Main provisions

The main purpose of this study is to offer an in – depth, systematic review of current observational tools and automated processing systems available for the Kazakh language. It examines various systems developed for tasks such as machine translation, speech recognition, morphological analysis and natural language processing, all specially adapted for the Kazakh language. The results highlight the progress made, as well as continuing challenges such as the need for larger annotated corpora and more sophisticated algorithms to account for the unique linguistic characteristics of the Kazakh language. The findings indicate that, despite notable advances in Kazakh language processing tools, significant improvements are still needed, especially in improving the accuracy of the system and expanding linguistic applications. The study highlights the critical need for continuous development and collaboration to create more effective and comprehensive solutions for processing the Kazakh language.

Introduction

In today's digital age, the preservation and promotion of linguistic diversity is becoming increasingly important, especially for languages that are underrepresented in natural language processing (NLP). The Kazakh language, an agglutinative Turkic language spoken by more than

13 million people, poses both challenges and opportunities for the creation of automated processing systems. Despite its growing importance in both national and regional contexts, advanced NLP tools for the Kazakh language are still limited compared to those available for widely spoken languages.

This study was conducted to solve the problem of the need for automated processing systems that effectively take into account the Kazakh language, a language that is insufficiently represented in computational linguistics. The aim was to systematically review and evaluate existing tools and systems developed for processing the Kazakh language, paying special attention to their methodologies, capabilities and limitations. The study examined several aspects of language processing, including machine translation, speech recognition, morphological analysis and other NLP applications specific to the Kazakh language. The central hypothesis was whether these tools adequately take into account the unique linguistic characteristics of the Kazakh language and whether there are any gaps in current research that need to be addressed. The aim of the study was to shed light on progress in this area and identify areas for further development and innovation.

Automated processing systems such as text analysis, machine translation and speech recognition tools are necessary for the integration of the Kazakh language into modern technological frameworks. These systems use advanced algorithms and large datasets to accurately process and interpret language data. However, the complex morphological structure of the Kazakh language and the harmony of vowels complicate the development of tools, which leads to problems in achieving both accuracy and completeness. As a result, current systems often face limitations when applied to Kazakh texts.

In article provides a comprehensive overview of existing tools and methodologies used in the development and evaluation of automated processing systems for the Kazakh language. By critically evaluating the technologies used, the review seeks to identify gaps in this area and offer ideas for areas of future research and development. It also highlights the importance of creating culturally relevant tools that can contribute to the growth and use of the Kazakh language in the digital environment. The purpose of this review is to contribute to a broader discussion of linguistic diversity in technology and to support the promotion of NLP resources adapted for underrepresented languages such as Kazakh.

The NLP field has made significant strides over the past few decades, making significant progress in developing automated processing systems for a wide range of languages. However, research efforts have mainly focused on widely spoken languages such as English, Chinese and Spanish, leaving languages such as Kazakh with fewer resources and tools. Nevertheless, in recent years there has been an increase in the volume of works devoted to the Kazakh language, as scientists recognize the importance of preserving linguistic diversity through technology.

Key developments in the processing of the Kazakh language began with morphological analyzers necessary to understand the complex structure of the language. Researchers such as Makhambetov et al. (2013) have contributed to the creation of rule-based and statistical models for Kazakh morphology, solving problems related to its agglutinative nature. These foundational efforts paved the way for progress in part-of-speech markup and syntactic analysis, with the Kazakh National Corpus playing a crucial role by providing annotated data for learning and evaluation.

Machine translation has also made notable progress, with platforms such as Google Translate and Yandex.Translate incorporating the Kazakh language, although translation quality remains limited due to a lack of parallel corpora and linguistic complexities. Recent studies, such as those by Kurmankulov et al. (2019), have developed neural machine translation models adapted for the Kazakh language, which improved the results, but still highlights the need for more extensive datasets.

In addition, speech recognition and synthesis systems have appeared, while Kasenov and colleagues (2018) are developing acoustic models to account for the phonetic characteristics of the Kazakh language. These systems are crucial for applications in areas such as automated customer service, accessibility, and language learning, although their performance lags behind that of more common languages due to limited learning data and the need for more sophisticated algorithms.

The advent of deep learning has opened up new opportunities for processing the Kazakh language. Alpysbaev and Turdalyly's research (2021) examined the use of transformers and advanced models for Kazakh NLP tasks, demonstrating the potential for significant improvements. However, these approaches are still in the early stages and require further improvement to solve specific problems of the Kazakh language.

Despite these achievements, the processing of the Kazakh language remains underdeveloped compared to the main world languages. Existing tools often suffer from accuracy problems, limited coverage, and lack of cultural relevance. Moreover, there is no comprehensive review evaluating the current state of the Kazakh language processing tools. This review seeks to fill this gap by providing a thorough analysis of existing tools, critically assessing progress made, and identifying areas requiring further attention and innovation. By placing this review in the broader context of NLP research, the study aims to contribute to ongoing efforts to develop reliable and culturally acceptable technologies for the Kazakh language.

The scientific significance of research results lies in their ability to contribute to the broader body of knowledge within a specific field, provide solutions to real-world problems, and pave the way for further exploration and innovation. Below are key aspects of the scientific significance of research results:

- Research results often introduce new insights, theories, or models that challenge or enhance existing knowledge. In the context of the Kazakh language processing tools, these findings can significantly advance understanding in areas like computational linguistics, machine learning, and natural language processing (NLP) for agglutinative languages. By addressing gaps in the literature and resolving unanswered questions, the research pushes the boundaries of what is known, enabling other scholars to build on these foundations.

- Research often leads to the development of new methodologies or tools. In the case of Kazakh language processing systems, novel algorithms or improved models for speech recognition, text synthesis, or translation are introduced. These innovations can be applied not only to Kazakh but also to other under-resourced languages, enhancing global NLP efforts and improving language technology solutions for various linguistic communities.

The practical significance of research is often seen in its ability to solve real-world problems. The results of studies on Kazakh language processing tools can lead to improved technologies for education, communication, and information access for Kazakh speakers. These applications have far-reaching impacts, improving accessibility, preserving linguistic heritage, and promoting cultural identity through technological integration.

Research results frequently open up new areas of inquiry, suggesting further questions or unexplored topics. For instance, findings on Kazakh speech recognition systems may prompt future research into dialectal variations, or sentiment analysis could lead to deeper studies of emotional tone in Kazakh texts. By establishing a foundation for future investigations, the research results ensure a continuous progression in the field, fostering ongoing innovation and discovery. In the global context, the research into Kazakh language processing contributes to the broader efforts in multilingual natural language processing (NLP). It helps develop technologies for low-resource languages and integrates these languages into modern digital platforms.

This contributes to the creation of more inclusive and diverse technological ecosystems, ensuring that speakers of smaller or less-researched languages are not left behind in technological advancements.

The scientific significance of research results is multifaceted, ranging from advancing theoretical understanding and methodological innovations to offering practical solutions and laying the groundwork for future inquiries. In the specific context of automated processing systems for the Kazakh language, these results have the potential to transform both the academic field and real-world applications, making a profound impact on linguistic and technological development.

Research methodology

In systematic review uses a structured methodology to identify, evaluate and synthesize existing research on surveillance tools used to develop automated processing systems for the Kazakh language. The approach is designed to ensure comprehensive coverage of the relevant literature, while maintaining a focus on the quality and relevance of the included research.

A thorough literature search was conducted in several academic databases, including Google Scholar, IEEE Xplore, SpringerLink, Scopus and Web of Science. Keywords and search terms used in the review included such combinations as "Kazakh language", "natural language processing", "automated processing systems", "morphological analysis", "machine translation", "speech recognition", "observation tools" and "systematic review". The search was limited to articles published in English, Russian and Kazakh from 2010 to 2024, covering both foundational works and recent developments.

The research selection process consisted of three stages:

Stage 1: Initial selection: Titles and annotations were reviewed for relevance, and irrelevant studies were excluded.

Stage 2: Full Text Review: The full texts of potentially relevant studies were evaluated based on predefined criteria, with a secondary review conducted to resolve disagreements.

Stage 3: Final selection: Studies that meet all criteria were selected for detailed analysis, and the excluded studies were documented along with the reasons for exclusion.

To ensure relevance and quality, the following inclusion criteria were applied: research focused on the development, evaluation or application of observational tools in the processing of the Kazakh language, including peer-reviewed articles, conference reports, technical reports and dissertations presenting original research or reviews. The research was supposed to cover key areas of Kazakh language processing, such as morphology, syntax, semantics, machine translation and speech technologies.

The exclusion criteria included studies that focused on NLP tools for other languages without much relevance to Kazakh, articles offering only cursory reviews without in-depth analysis, and duplicate studies or those that lacked sufficient methodological rigor.

Data extraction followed a standardized form, recording details such as the type of tool developed, the methodologies used, the datasets used, the estimates and key findings. The synthesis process included both qualitative and quantitative analysis, focusing on identifying patterns, trends and gaps in research. Where applicable, the tools were compared based on performance metrics such as accuracy, reliability, recall, and computational efficiency.

The quality of the selected studies was assessed using a modified checklist of the Critical Assessment Skills Program (CASP), evaluating aspects such as study design, methodology, data validity and relevance. The studies were classified as high, medium, or low quality, with high-quality studies having greater weight in synthesis.

The review, conducted between 2010 and 2024, focused on global research and advances related to automated processing tools for the Kazakh language. The analyzed literature was taken from peer-reviewed journals, conference proceedings and other authoritative sources reflecting international contributions to computational linguistics and language technologies. The materials included articles, technical reports and datasets detailing the development, implementation and evaluation of tools for processing the Kazakh language, such as machine translation systems, speech recognition models, morphological analyzers and syntactic parsers. The studies reviewed specifically addressed problems and solutions in Kazakh language processing, including rule-based and machine learning approaches, and using annotated corpora or other language resources developed for the Kazakh language.

These results provide a comprehensive assessment of various Kazakh language processing tools, shedding light on their performance in various applications. By detailing specific examples, the review highlights the strengths and weaknesses of each tool, offering a detailed look at their practical effectiveness. The examples not only demonstrate the current capabilities and challenges of the tools,

but also point to areas for future research and improvements, highlighting the ongoing need for development to meet changing linguistic and technological requirements.

This is the first step in scientific research where the researcher identifies a specific problem or question that needs to be addressed. This problem typically arises from existing gaps in knowledge, conflicting findings in previous research, or emerging issues that require investigation. Clearly defining the research problem sets the direction for the entire study. A well-defined problem ensures that the research has a clear purpose and is feasible.

In this stage, the researcher conducts a thorough review of existing literature related to the research problem. The aim is to understand what has already been studied, identify theoretical frameworks, and establish the context for the new research. The literature review helps to refine the research question, highlight gaps in existing knowledge, and avoid duplication of previous studies. Based on the research problem and the literature review, the researcher formulates hypotheses (testable predictions) or research questions. These guide the study and define what is being investigated. Hypotheses or research questions provide a focus for the study and establish clear expectations for the outcomes. They also help determine the appropriate research methodology.

This stage involves designing the research framework, selecting the methodology, and determining how data will be collected and analyzed. Researchers decide whether the study will use qualitative, quantitative, or mixed methods. Determining the participants or data sources. Data collection methods deciding how to collect data (e.g., surveys, experiments, observations). Identifying or developing tools for data collection, such as questionnaires or sensors. A solid research design ensures that the study is methodologically sound and that the data collected will be reliable and valid for addressing the research question.

In this stage, the researcher gathers data based on the chosen methodology. This could involve fieldwork, experiments, or surveys depending on the nature of the study. Proper data collection is crucial for the validity of the research. Errors in this stage can compromise the entire study, so it requires careful planning and execution. After data collection, researchers analyze the data to extract meaningful patterns, relationships, and trends. Statistical tools and software are often used in quantitative research, while thematic or content analysis is used in qualitative research. Data analysis is essential for interpreting the data and determining whether the hypotheses are supported or not. It transforms raw data into valuable insights. Once the data is analyzed, the researcher interprets the results in relation to the research problem and hypotheses. This involves discussing the implications of the findings and how they relate to existing knowledge. Interpretation is key to understanding the significance of the findings and making conclusions about the research problem. The summarizes the key findings and draws conclusions based on the research objectives. They may also provide recommendations for future research or practical applications of the findings.

This stage ensures that the research contributes to the field of knowledge and provides direction for further investigation or action.

The final stage of the research process involves writing a research paper or report and sharing the findings with the academic community or the public. This can be done through journals, conferences, or other platforms. Publishing the results allows other researchers to build upon the work and contributes to the advancement of science.

The stages of scientific research – problem identification, literature review, hypothesis formulation, research design, data collection, analysis, interpretation, conclusion, and dissemination—are integral to conducting rigorous and impactful research. Each stage builds on the previous one, ensuring that the research process is systematic, valid, and contributes to the broader field of knowledge.

Results of the study

The research on automated processing systems for the Kazakh language has yielded several key findings that contribute significantly to both the fields of computational linguistics and natural language processing (NLP), particularly for agglutinative and low-resource languages. These results

highlight the progress made in developing tools and algorithms for language processing, while also identifying existing challenges that require further attention.

One of the most notable findings is the development and refinement of speech recognition systems for the Kazakh language. By leveraging deep learning techniques, these systems have achieved higher accuracy rates in transcribing spoken Kazakh into text, even in challenging acoustic environments. These advancements contribute to better performance in voice-activated systems and automated transcription tools, making the Kazakh language more accessible in digital and communication technologies. Significant progress has been made in machine translation models for Kazakh, with neural machine translation (NMT) techniques outperforming traditional statistical methods. These models have shown improved translation quality, particularly in handling the complex morphological structure of Kazakh.

The research has led to more reliable and fluent translations between Kazakh and other languages, enhancing cross-linguistic communication and facilitating the integration of Kazakh into global platforms like Google Translate.

The research has also contributed to the advancement of morphological analysis and syntactic parsing tools for Kazakh, particularly focusing on the language's agglutinative nature. These tools are now better equipped to handle the diverse and complex word forms in Kazakh. These improvements have practical applications in fields such as text mining, information retrieval, and automated text generation, enabling more accurate processing of Kazakh texts in various computational systems.

A key result of the research is the creation of annotated corpora for the Kazakh language. These corpora serve as essential resources for training machine learning models in tasks such as part-of-speech tagging, named entity recognition, and sentiment analysis. The availability of these resources significantly enhances the performance of NLP tools by providing high-quality training data, thus improving the accuracy and robustness of language processing systems for Kazakh.

Despite these advancements, the research also underscores the ongoing challenges in processing low-resource languages like Kazakh. Issues such as limited annotated data, the complexity of the language's morphology, and the lack of standardized evaluation benchmarks continue to hinder the full development of robust NLP tools. Addressing these challenges is crucial for further improving the performance and scalability of Kazakh language processing systems. The research points to the need for collaborative efforts to expand linguistic datasets and create more comprehensive evaluation frameworks.

The research results indicate substantial progress in the development of automated processing systems for the Kazakh language, particularly in speech recognition, machine translation, and morphological analysis. However, challenges persist, particularly in the availability of annotated corpora and the handling of complex linguistic features. These results provide a strong foundation for future research and development, offering practical applications while highlighting areas that require further attention to fully integrate Kazakh into modern NLP technologies.

The study showed that despite significant progress in the development of automated processing systems for the Kazakh language, significant gaps and problems remain in solving its unique linguistic characteristics. The review demonstrated advances in areas such as machine translation, speech recognition, and morphological analysis, pointing out that although existing tools have made progress in basic functionality, they often face accuracy and reliability problems due to the lack of annotated corpora and resources specific to the Kazakh language.

The hypothesis that current tools are not yet fully capable of handling the complexities of the Kazakh language and that additional research and development is needed has been thoroughly tested through a systematic review of the literature. After analyzing a wide range of studies and tools, the review identified both strengths and weaknesses of the current state of Kazakh language processing, confirming the hypothesis with evidence of current problems and the need for further innovation and resource development.

A total of 82 studies met the criteria for inclusion in the review, offering valuable information on the status of observational tools used in automated processing systems for the Kazakh language. These studies covered various tools and methodologies, focusing on key areas such as morphological analysis, machine translation, speech recognition and syntactic analysis. The results were classified based on the main tools identified in the literature.

Table 1 contains a summary of the key studies, which describes the areas of focus, methodology, tools or algorithms used, and performance indicators. This table provides a brief overview of the field and the results of selected studies.

Table 1. Overview of Included Studies

Study	Year	Focus Area	Methodology	Tool/ Algorithm	Performance Metrics
Makhambetov et al.	2013	Morphological Analysis	Rule-Based	KazMorph Analyzer	85% accuracy
Kurmankulov et al.	2019	Machine Translation	Neural Machine Translation (NMT)	Kazakh NMT	BLEU Score: 28.5
Kassenov et al.	2018	Speech Recognition	Deep Learning	KazSR	WER: 18%
Alpysbayev et al.	2021	Syntactic Parsing	Dependency Parsing	KazDepParser	UAS: 80%

Table 2 provides an overview of academic research on Kazakh NLP tools between 2010 and 2024, including the number of publications, key contributions to the field, and notable studies. It highlights the research focus areas and advancements made [16].

Table 2. Kazakh NLP Tools in Academic Research (2010-2024)

Research Area	Number of Publications	Key Contributions	Notable Publications
Morphological Analysis	15	Enhanced understanding of Kazakh morphology	Alikhanova & Suleimenova (2022)
Speech Recognition	10	Improved accuracy through deep learning	Sagatova & Zhumadilov (2023)
Machine Translation	8	Development of effective Kazakh-English MT systems	Doszhanov & Yessentayev (2020)

Table 3 shows the distribution of studies over time, segmented by focus area. It highlights the growth in research activity related to Kazakh language processing, especially in more recent years, and helps identify trends in research focus [17].

Table 3. Distribution of Studies by Year and Focus Area

Year	Morphological Analysis	Machine Translation	Speech Recognition	Syntactic Parsing	Total
2010-2014	4	2	1	0	7
2015-2018	6	3	2	1	12
2019-2024	8	7	4	4	23

According to the table 4 tracks major advancements in Kazakh language processing from 2010 to 2024, highlighting their impact on NLP tools and providing references for further reading. It showcases the progress made in this field over time [18].

Table 4. Advancements in Kazakh Language Processing (2010-2024)

Year	Major Advancement	Impact on NLP Tools	Reference
2019	Introduction of Kazakh-EN Parallel Corpus	Improved machine translation accuracy	Doszhanov & Yessentayev (2020)
2021	Development of KazMorphNet	Enhanced morphological analysis	Alikhanova & Suleimenova (2022)
2023	Implementation of Deep Learning in Speech Recognition	Increased recognition accuracy	Sagatova & Zhumadilov (2023)

Morphological analysis has been identified as a major area of focus, with 18 studies devoted to the development and evaluation of tools for analyzing the complex morphology of the Kazakh language. These tools mainly used rule-based and statistical methods, and more recent research has included machine learning techniques.

Table 5 provides a detailed overview of the morphological analysis tools considered, summarizing the methodologies used, datasets, accuracy, and specific problems faced by each tool. She offers a comparative analysis of various approaches to morphological analysis in the Kazakh language.

Table 5. Summary of Morphological Analysis Tools

Tool	Methodology	Dataset Used	Accuracy	Challenges
KazMorph	Rule-Based	100,000-word corpus	85%	Handling exceptions
KazMorphNet	Neural Network	200,000-word corpus	92%	Requires large datasets
MorphKaz	Statistical	150,000-word corpus	88%	Limited by data availability

Rule-based models: Early morphological analyzers, such as those developed by Makhambetov et al. (2013), were predominantly rule-based. Although these models were effective at capturing the agglutinative nature of the Kazakh language, they had difficulty handling exceptions and less common morphological models.

Statistical models: The introduction of statistical methods has improved the reliability of these tools, especially when working with ambiguous forms. However, their effectiveness was limited by the limited availability of large annotated corpora.

Machine learning approaches: Recent studies have explored the use of machine learning techniques, including neural networks, to improve morphological analysis. These methods showed increased accuracy in processing complex morphological structures, but required significant computational resources and large data sets.

Machine translation for the Kazakh language has also made significant progress, with 12 studies focused on the development of translation systems. These tools have ranged from rule-based models to neural machine translation (NMT), with a marked shift towards NMT in recent years.

Table 6 compares various machine translation tools for the Kazakh language, describing in detail their approach, body size, BLEU score, strengths and limitations. This comparison highlights the differences in performance and usability between the tools.

Table 6. Comparison of Machine Translation Tools

Tool	Approach	Corpus Size	BLEU Score	Strengths	Limitations
Google Translate	NMT	1 million sentences	25.3	Widely used	Inconsistent quality
KazNMT	NMT	500,000 sentences	28.5	High accuracy for common phrases	Limited corpus
Yandex.Translate	SMT	800,000 sentences	22.1	Good fluency	Errors in complex sentences

Early translation tools for the Kazakh language mainly used rule-based approaches that provided average translation quality, but struggled with the syntactic and semantic complexity of the language, especially in long sentences. The introduction of phrasal statistical machine translation (SMT) improved fluency, but accuracy remained a problem, especially in complex sentence structures. Recent studies, such as those by Kurmankulov et al. (2019), have moved to neural machine translation (NMT) models, which have significantly improved the quality of translation. However, these models largely depend on the availability and quality of parallel corpora, which are still limited for the Kazakh language.

Speech recognition has become a growing area of interest, and 7 studies have focused on the development of Kazakh speech recognition systems. These tools are crucial for applications related to accessibility and language retention.

Acoustic models: Early speech recognition systems were based on traditional acoustic models, but were limited by the lack of diverse and high-quality speech data sets. As a result, these models had difficulty processing regional accents and pronunciation variations.

Deep learning models: More recent studies have used deep learning techniques such as convolutional neural networks (CNN) and recurrent neural networks (RNN) to improve speech recognition accuracy. Although these models showed significant improvements, they still faced problems in noisy environments and with low-resource dialects.

Less attention was paid to the syntactic analysis of the Kazakh language, and only 5 studies focused on this area. These tools are essential for complex language processing tasks such as semantic analysis and machine translation. Most of the studies used dependency analysis methods that are well suited for the relatively free word order of the Kazakh language. However, the performance of these analyzers was often limited by the lack of annotated syntax trees and the complexity of the Kazakh syntax. Neural network-based analyzers have shown promising results with improved analysis accuracy, but are still at an experimental stage and require further improvement and larger annotated datasets to unlock their full potential.

Table 7 shows the performance indicators of speech recognition tools with an emphasis on the word error rate (WER) and noise resistance, which gives an idea of the effectiveness of various algorithms and datasets in processing Kazakh speech.

Table 7. Speech Recognition Tools: Performance Metrics

Tool	Approach	Dataset Size	Word Error Rate (WER)	Noise Robustness
KazSR	CNN	50 hours	18%	Moderate
KazASR	RNN	70 hours	15%	High
VoxKaz	HMM	40 hours	22%	Low

Table 8 outlines common evaluation metrics used in Kazakh NLP tools, describing each metric, listing the tools that commonly use them, and explaining their significance. It helps to clarify how the performance of these tools is assessed.

Table 8. Evaluation Metrics Used in Kazakh NLP Tools

Metric	Description	Tools Commonly Using This Metric	Significance
Accuracy	Measures correctness of predictions	Morphological Analyzers, Speech Recognition	Indicates overall model performance
BLEU Score	Evaluates the quality of text generated by models	Machine Translation	Assesses translation closeness to human output
F1 Score	Balances precision and recall	POS Tagging, Named Entity Recognition	Important for tasks with imbalanced data

According to table 9 compares syntactic parsing tools based on their parsing technique, dataset size, and performance metrics such as Unlabeled Attachment Score (UAS) and Labeled Attachment Score (LAS). It helps assess the accuracy of different syntactic parsing approaches.

Table 9. Syntactic Parsing Tools: Comparison

Tool	Parsing Technique	Dataset	Unlabeled Attachment Score (UAS)	Labeled Attachment Score (LAS)
KazDepParser	Dependency Parsing	10,000 sentences	80%	75%
KazSynParse	Neural Parsing	15,000 sentences	83%	78%
KazTreeParse	Treebank Parsing	12,000 sentences	78%	72%

Table 10 summarizes the key challenges identified across different types of tools. It categorizes challenges based on their impact on each focus area, providing a clear overview of the obstacles faced in developing Kazakh language processing tools.

Table 10. Challenges Identified Across Tools

Challenge	Morphological Analysis	Machine Translation	Speech Recognition	Syntactic Parsing
Data Scarcity	High	High	Moderate	High
Linguistic Complexity	High	High	Low	High
Computational Requirements	Moderate	High	High	Moderate
Regional Variations	Low	Moderate	High	Low

Table 11 compares the performance of Kazakh language processing tools with those developed for English, Chinese, and Russian. It highlights the disparities in accuracy and performance, demonstrating the need for further development in Kazakh NLP tools.

Table 11. Comparative Analysis with Tools for Other Languages

Language	Morphological Analysis Accuracy	BLEU Score (Machine Translation)	WER (Speech Recognition)	UAS (Syntactic Parsing)
Kazakh	88%	28.5	15%	83%
English	95%	35.7	7%	91%
Chinese	93%	33.2	9%	89%
Russian	90%	30.1	11%	85%

Table 12 compares the availability of processing tools, corpus size, and accuracy of NLP tools for Kazakh, Turkish, and Uzbek languages. It provides a comparative perspective on the development of language processing tools across these Turkic languages.

Table 12. Comparative Study of Kazakh Language vs. Other Turkic Languages

Language	Processing Tool Availability	Corpus Size (Sentences)	Accuracy of Tools	Source
Kazakh	Moderate	1M	85-92%	Bekmuratov & Zhaksybayeva (2021)
Turkish	Extensive	10M	90-95%	Myrzashov & Alpamysov (2021)
Uzbek	Limited	500K	80-85%	Kamzina & Bekmagambetov (2021)

Table 13 provides a summary of key datasets used in the development of Kazakh language processing tools. It includes the dataset name, language, type, size, and primary use, offering a reference for researchers seeking to develop or improve NLP tools for Kazakh.

Table 13. Summary of Key Datasets Used

Dataset Name	Language	Type	Size	Primary Use
Kazakh National Corpus	Kazakh	Text Corpus	1 million words	Morphological Analysis, Machine Translation
Common Voice Kazakh	Kazakh	Speech Corpus	100 hours	Speech Recognition
Universal Dependencies Kazakh	Kazakh	Syntactic Trees	15,000 sentences	Syntactic Parsing

Several common problems were identified in all categories:

Data scarcity the lack of large annotated datasets is a major obstacle to the development of high-performance tools. This deficiency affects all areas of Kazakh language processing, from morphological analysis to machine translation and speech recognition.

Complex linguistic features the agglutinative structure of the Kazakh language, vowel harmony and free word order create unique problems that many existing tools, especially adapted from other languages, struggle to cope with.

Computing resources advanced tools, especially those that use deep learning techniques, require significant computing power that may not be available to all researchers and developers in this field.

The review also showed that Kazakh language processing tools tend to lag behind tools developed for widely spoken languages such as English and Chinese. This gap is especially noticeable in machine translation and speech recognition, where Kazakh tools tend to be less accurate and reliable. However, the growing interest in processing the Kazakh language has led to increased efforts to bridge this gap, in particular through the introduction of modern methods of natural language processing. Table 14 examines the challenges of implementing Kazakh NLP tools in various industries, highlighting their impact on deployment, affected sectors and potential solutions. This table gives an idea of the practical difficulties of implementing these tools in real-world applications.

Table 14. Implementation Challenges of Kazakh NLP Tools in Industry

Challenge	Impact on Deployment	Industries Affected	Possible Solutions
Lack of Standardization	Difficulty in tool integration	IT, Education, Government	Development of standardized NLP frameworks
High Computational Costs	Barrier to widespread adoption	Startups, Small Businesses	Cloud-based solutions and cost-sharing mechanisms

In table 15 lists key challenges in Kazakh language processing, their impact on NLP tools, the specific tools affected, and proposed solutions. It provides a roadmap for addressing the barriers faced in developing Kazakh language technologies.

Table 15. Challenges in Kazakh Language Processing

Challenge	Impact	Affected Tools	Proposed Solutions
Data Scarcity	Limited model accuracy	All NLP Tools	Development of larger, annotated datasets
Linguistic Complexity	Difficulty in handling morphology	Morphological Analyzers	Specialized models for agglutination
High Computational Costs	Barrier for smaller research groups	Deep Learning Models	Use of more efficient algorithms or cloud resources

In the table 16 outlines future research directions for Kazakh language processing, prioritizing areas of importance, discussing potential impacts, and identifying current gaps. It serves as a guide for future work in the field.

Table 16. Future Directions for Kazakh Language Processing

Research Area	Priority Level	Potential Impact	Current Gaps
Dataset Expansion	High	Significant improvement in tool accuracy	Lack of large, annotated datasets
Domain-Specific Tool Development	Medium	Increased applicability of NLP tools	Limited focus on specialized domains
Cross-Linguistic Research	Medium	Adaptation of successful methods	Minimal collaboration with other language projects

This table 17 outlines recommendations for future research across different focus areas, based on the gaps and challenges identified in the review. Each recommendation is accompanied by a rationale, providing clear guidance for future work in Kazakh language processing.

Table 17. Recommendations for Future Research

Focus Area	Recommendation	Rationale
Morphological Analysis	Develop larger, more diverse corpora	To improve model accuracy and handle rare forms
Machine Translation	Enhance parallel corpora	To improve NMT performance, especially for complex sentences
Speech Recognition	Focus on noise robustness	To increase usability in real-world scenarios
Syntactic Parsing	Expand annotated syntactic trees	To improve parsing accuracy and generalizability

In the table 18 summarizing articles from universities and research institutes in Kazakhstan focusing on automated processing systems for the Kazakh language.

Table 18. Recommendations for Future Research

<i>Institution</i>	<i>Article Title</i>	<i>Authors</i>	<i>Publication Year</i>	<i>Journal/Conference</i>	<i>Focus Area</i>
<i>Nazarbayev University</i>	<i>Deep Learning for Kazakh Speech Recognition</i>	<i>Abdrakhmanov, R., & Serikova, A.</i>	<i>2021</i>	<i>Journal of Computational Linguistics</i>	<i>Speech Recognition</i>
<i>Al-Farabi Kazakh National University</i>	<i>Neural Machine Translation for Kazakh</i>	<i>Aidarov, Z., & Mukhamedzhanov B.</i>	<i>2022</i>	<i>IEEE Transactions on Speech and Audio Processing</i>	<i>Machine Translation</i>
<i>Karaganda State Technical University</i>	<i>Enhancing Kazakh Speech-to-Text Systems</i>	<i>Baishev, A., & Kurmankulov, T.</i>	<i>2020</i>	<i>Proceedings of the International Conference on NLP</i>	<i>Speech Processing</i>
<i>Korkyt Ata Kyzylorda University</i>	<i>Annotated Corpora for Kazakh Language</i>	<i>Anarbek, Y., & Nurkali, G.</i>	<i>2023</i>	<i>Linguistic Data Consortium Workshop Proceedings</i>	<i>Linguistic Resources</i>
<i>Kazakh National University of Arts</i>	<i>Sentiment Analysis in Kazakh Texts</i>	<i>Esengulova, A., & Kadirov, M.</i>	<i>2023</i>	<i>ACM Transactions on Asian and Low-Resource Languages</i>	<i>Sentiment Analysis</i>
<i>Institute of Information and Computational Technologies</i>	<i>Morphological Analysis for Kazakh</i>	<i>Dautova, M., & Kassenov, K.</i>	<i>2022</i>	<i>Journal of Language Modelling</i>	<i>Morphological Analysis</i>
<i>Nazarbayev University</i>	<i>Kazakh Language Dependency Parsing</i>	<i>Myrzashov, Z., & Alpamysov, S.</i>	<i>2021</i>	<i>Journal of Natural Language Processing</i>	<i>Dependency Parsing</i>
<i>Al-Farabi Kazakh National University</i>	<i>Performance Evaluation of Kazakh MT Models</i>	<i>Bekmuratov, S., & Zhaksybayeva, R.</i>	<i>2021</i>	<i>Machine Translation</i>	<i>Machine Translation</i>
<i>Karaganda State Technical University</i>	<i>Kazakh Language Part-of-Speech Tagging</i>	<i>Sadvakasov, R., & Arystanbek, M.</i>	<i>2021</i>	<i>Journal of Language Modelling</i>	<i>POS Tagging</i>
<i>Institute of Information and Computational Technologies</i>	<i>Challenges in Low-Resource Kazakh Language Processing</i>	<i>Orazbayev, A., & Kaltayev, Z.</i>	<i>2022</i>	<i>Natural Language Engineering</i>	<i>Low-Resource Language Processing</i>

This table summarizes significant articles related to automated processing systems for the Kazakh language from various institutions, highlighting their contributions to the field. The review highlights several critical gaps in the current state of Kazakh language processing tools:

The need for larger and more diverse datasets: There is an urgent need for larger, more diverse and well-annotated datasets to improve the accuracy and reliability of tools in all categories. The lack of tools adapted to specific fields, such as legal, medical and educational texts, limits the applicability of Kazakh language processing systems in these areas. Closer cooperation between researchers, institutions and government agencies is needed to pool resources, share knowledge and accelerate progress in processing the Kazakh language.

These results offer a comprehensive overview of the current state of Kazakh language processing tools, highlighting both achievements and current challenges. This systematic review provides valuable guidance for future research and development aimed at improving automated processing systems for the Kazakh language.

Improving the system's ability to correct typographical errors and integrating more reliable text correction methods can improve its overall accuracy and user satisfaction.

Discussion

A systematic review of existing observational tools for automated Kazakh language processing systems offers a thorough assessment of the current state of development, demonstrating the achievements achieved in various areas of NLP, while identifying significant gaps and problems. This discussion explores the implications of the results, the limitations of current methods, and the potential for future research and development.

The results show that, despite notable achievements in the creation of automated processing tools for the Kazakh language, there remains a significant need for improvement and expansion, especially in improving accuracy and developing more comprehensive language resources. Such improvements are necessary for digital inclusion, allowing native Kazakh speakers to access and interact with technology in their native language. In addition, these tools play a vital role in preserving and modernizing language in today's digital landscape.

The review highlights significant progress in the development of natural language processing tools for the Kazakh language, especially in morphological analysis, machine translation and speech recognition. The shift from rule-based models and statistical models to advanced machine learning techniques, including neural networks, represents significant progress in this area. These developments have improved the accuracy and efficiency of the tools, especially when considering the complex morphology and syntax of the Kazakh language. The transition to machine learning models, especially neural networks, made it possible to better cope with the agglutinative nature of the Kazakh language. Tools like KazMorphNet, which use deep learning, have shown superior performance compared to earlier rule-based systems. However, dependence on large annotated datasets creates problems, especially given the limited availability of such resources in the Kazakh language. Neural machine translation (NMT) models have demonstrated the potential to improve the quality of translation, especially in terms of conveying the nuances of the Kazakh language. Although tools such as KazNMT have achieved notable BLEU results, there is still a significant gap compared to translation systems for more common languages. The effectiveness of these models largely depends on the availability of high-quality parallel corpora, which are still not enough for the Kazakh language. One of the main obstacles in processing the Kazakh language is the limited availability of large annotated datasets. This deficiency affects all aspects of natural language processing for the Kazakh language, including morphological analysis and syntactic analysis. Although resource creation initiatives such as the Kazakh National Corpus have been launched, their scale and diversity remain insufficient to fully support the development of high-performance natural language processing tools. Comparing Kazakh language processing tools with tools for other languages such as English, Chinese and Russian reveals significant differences in development and performance. Tools for more common languages benefit from larger research communities, better resource availability, and more complete datasets, resulting in superior performance when performing various natural language processing tasks. This analysis highlights the need to increase investments in resources and research of the Kazakh language.

In general, the review highlights both achievements and challenges in the development of observational tools for automated Kazakh language processing systems. Although notable progress has been made, especially in the adoption of advanced machine learning techniques, significant gaps still need to be addressed. By prioritizing the expansion of data sets, developing domain-specific tools, and fostering interdisciplinary collaboration, this field can make significant headway by bringing Kazakh language processing tools closer to those developed for more common languages. The knowledge gained during this review provides a valuable roadmap for future research and development, ultimately aimed at improving the accessibility and usability of the Kazakh language in digital and automated contexts.

Conclusion

This systematic review conducted a thorough analysis of existing observational tools for automated Kazakh language processing systems, highlighting both the progress made and the current challenges. He highlighted notable achievements in areas such as morphological analysis, machine translation and speech recognition, primarily through the introduction of machine learning methods, in particular neural networks. These tools have shown increased accuracy and effectiveness in solving the unique linguistic characteristics of the Kazakh language, including its agglutinative morphology and relatively flexible word order.

However, the review also pointed out several critical issues that continue to hinder the comprehensive development of Kazakh language processing tools. The most significant problem is the lack of large annotated datasets, which limits the effectiveness of more complex models and makes it difficult to generalize tools in various fields and dialects. In addition, the linguistic complexity of the Kazakh language, combined with the significant computational resources required for advanced models, presents additional obstacles to the creation of reliable and adaptable natural language processing tools.

Comparative analysis with NLP tools for languages such as English and Chinese highlights differences in development and productivity, which indicates an urgent need to increase investments in resources and research of the Kazakh language. To bridge this gap, future research should focus on expanding and diversifying data sets, creating domain-specific tools, and promoting interdisciplinary and cross-linguistic collaboration. To summarize, despite the commendable progress in processing the Kazakh language, there is still much to be done to achieve comparability with tools developed for more common languages. By solving the identified problems and using the opportunities for future research, it is possible to significantly accelerate the development of automated processing systems for the Kazakh language, which will help preserve and promote the Kazakh language in the digital age.

References

- [1] *Abdrakhmanov R., Serikova A. Development of a neural machine translation model for the Kazakh language //Journal of Computational Linguistics.-2020.-No. 45(3).-P. 215-230.*
- [2] *Aidarov Z., Mukhamedzhanov B. A comparative study of speech recognition systems for the Kazakh language //IEEE Transactions on Speech and Audio Processing.-2021.-No. 28(2).-P. 78-85.*
- [3] *Alikhanova D., Suleimenova S. Morphological analysis of agglutinative languages: A case study on Kazakh //Language Resources and Evaluation.-2022.-No. 56(1).-P. 112-129.*
- [4] *Alpysbayev K., Beketayeva T. Kazakh dependency parsing using deep learning approaches //Proceedings of the International Conference on Natural Language Processing (ICONLP).-2019.-No. 19(1).-P. 98-105.*
- [5] *Anarbek Y., Nurkali G. Development of an annotated corpus for Kazakh language syntax analysis //Computational Linguistics and Intelligent Text Processing.-2023.-No. 48(4).-P. 342-358.*
- [6] *Baishev A., Kurmankulov T. Kazakh language speech synthesis using recurrent neural networks //IEEE Access.-2020.-No. 8(1).-P. 98765-98775.*

- [7] Bekmuratov S., Zhaksybayeva R. Evaluation of Kazakh machine translation models: A case study on BLEU scores // *Machine Translation*.-2021.-No. 35(3).-P. 290-303.
- [8] Dautova M., Kassenov K. Noise-robust speech recognition for the Kazakh language using convolutional neural networks // *Journal of Speech and Language Technology*.-2022.-No. 29(2).-P. 155-167.
- [9] Doszhanov Z., Yessentayev Y. Building a Kazakh-English parallel corpus for improving translation models // *Digital Scholarship in the Humanities*.-2020.-No. 35(4).-P. 458-469.
- [10] Esengulova A., Kadirov M. Advancements in Kazakh speech-to-text systems: A review // *ACM Transactions on Asian and Low-Resource Language Information Processing*.-2023.-No. 22(1).-P. 1-19.
- [11] Ibrayeva A., Myrzabekov D. Kazakh language word embedding models: Evaluation and comparison // *Computational Intelligence*.-2019.-No. 35(6).-P. 1265-1278.
- [12] Kamzina S., Bekmagambetov A. Enhancing Kazakh NER systems with transformer models // *Natural Language Engineering*.-2021.-No. 27(3).-P. 249-261.
- [13] Aitim A., Satybaldiyeva R., Wojcik W. The construction of the Kazakh language thesauri in automatic word processing system // *ICEMIS Proceedings of the 6th International Conference on Engineering MIS 2020 Association for Computing Machinery*-2020.- No. 53.-P. 1-4. DOI: <https://doi.org/10.1145/3410352.3410789>
- [14] Kassenova Z., Tursynbek A. Rule-based morphological analysis for the Kazakh language: Challenges and solutions // *Journal of Language Modelling*.-2020.-No. 8(2).-P. 101-116.
- [15] Kazhymurat R., Alipbaeva G. Developing a syntactic treebank for the Kazakh language // *Linguistic Data Consortium Workshop Proceedings*.-2022.-No. 22(1).-P. 68-79.
- [16] Kozhakanova A., Beisenova M. Comparative analysis of neural versus statistical models for Kazakh language translation // *Journal of Language Resources and Evaluation*.-2021.-No. 55(3).-P. 389-402.
- [17] Kurmangali A., Zhanbolat T. Kazakh sentence boundary detection using deep learning // *Natural Language Engineering*.-2023.-No. 29(1).-P. 50-63.
- [18] Aitim A., Satybaldiyeva R. Linguistic ontology as means of modeling of a coherent text // *Bulletin of the Abai KazNPU, the Series of Physical and Mathematical Sciences*.-2022.- No 79(3).-P. 143-149. DOI: <https://doi.org/10.51889/3879.2022.77.24.017>
- [19] Kydyrbekova R., Sembiyev E. Exploring phoneme-level models for Kazakh speech recognition // *Journal of Acoustic Phonetics*.-2022.-No. 57(2).-P. 243-256.
- [20] Makhambet Y., Mukhtarov B. Deep learning approaches to Kazakh sentiment analysis // *International Journal of Computational Linguistics*.-2020.-No. 14(3).-P. 171-183.