

## ПРИМЕНЕНИЕ ГЛУБОКИХ НЕЙРОННЫХ СЕТЕЙ ДЛЯ ЛИНГВИСТИЧЕСКОЙ РЕКОНСТРУКЦИИ КАЗАХСКОГО ТЕКСТА

Омаргазы Данияр Ермакулы

магистрант, Казахстанско-Британский технический университет, Казахстан, г. Алматы

**Аннотация.** В данной статье рассматриваются возможности применения глубоких нейронных сетей для лингвистической реконструкции текстов на казахском языке. Описаны основные архитектуры и методы, позволяющие моделировать сложные лингвистические и культурные особенности казахского текста, включая грамматическую агглютинацию и исторические контексты. Особое внимание уделяется адаптации современных NLP-моделей к уникальным характеристикам казахского языка и процессам реконструкции, которые способствуют более глубокому пониманию и сохранению культурного наследия. Приведены примеры успешной интерпретации и преобразования текстов, а также оценка эффективности методов. Исследование подчеркивает значимость интеграции технологий глубокого обучения для анализа и возрождения казахского текстового наследия, открывая новые перспективы для его изучения и адаптации в цифровую эпоху.

**Ключевые слова:** глубокие нейронные сети, казахский язык, лингвистическая реконструкция, обработка естественного языка, культурное наследие, NLP, агглютинация, машинное обучение, реконструкция текста, искусственный интеллект.

Сохранение и понимание культурного наследия становятся все более актуальными задачами в условиях стремительного технологического развития и глобализации. Казахский язык, как и многие другие, нуждается в современных инструментах для анализа и восстановления культурных и исторических смыслов, заложенных в текстах прошлых эпох. Применение глубоких нейронных сетей открывает новые возможности для лингвистической реконструкции казахских текстов, позволяя выявлять и интерпретировать сложные лингвистические и культурные аспекты, которые могли бы быть утрачены или не поняты современными читателями [1].

Цель данной работы — исследовать возможности применения глубоких нейронных сетей для автоматизированного анализа, интерпретации и реконструкции казахских текстов. Специфика казахского языка, включая агглютинацию, богатую морфологию и особую синтаксическую структуру, требует адаптации современных моделей обработки естественного языка (NLP), что делает задачу как интересной, так и научно значимой.

Использование глубоких нейросетевых моделей для казахского языка позволяет не только восстанавливать смыслы и улучшать понимание текста, но и способствует развитию технологий сохранения языков и культурного наследия в условиях глобализации.

Лингвистическая реконструкция представляет собой процесс восстановления возможного оригинального значения, структуры или звучания текста, что особенно важно для языков, богатых культурными кодами [2]. В контексте казахского языка реконструкция может быть применена для адаптации исторических и классических текстов к современным условиям, что позволяет сохранить уникальные культурные аспекты, не утрачивая смысла.

Глубокие нейронные сети, в особенности модели обработки естественного языка (NLP), представляют собой перспективный инструмент для лингвистической реконструкции благодаря своей способности выявлять скрытые зависимости и сложные языковые структуры. Наиболее популярные архитектуры, такие как BERT, GPT и T5, успешно применяются для анализа текстов на различных языках, обеспечивая возможность не только их анализа, но и генерации текста с учетом заданных параметров. Для казахского языка, однако, необходима дополнительная адаптация моделей, чтобы учесть морфологические и синтаксические особенности, уникальные для этой языковой системы.

Специфика казахского языка в контексте анализа нейросетями

Казахский язык обладает особенностями, которые представляют определенные вызовы для NLP-моделей. К таким характеристикам относятся:

- **Агглютинация:** Казахский язык является агглютинативным, что означает присоединение морфем к основе слова для изменения его значения или функции. Это требует от моделей способности учитывать различные морфемные структуры, создаваемые путем сложных комбинаций основ и аффиксов.
- **Грамматическая сложность и многозначность:** Некоторые казахские слова могут менять значение в зависимости от контекста, и их перевод может требовать учета культурных особенностей, что усложняет автоматическую обработку текста.
- **Историческая вариативность:** В казахских текстах разных эпох могут встречаться устаревшие формы слов или выражения, требующие лингвистической адаптации для их понимания современными читателями.

Для эффективной реконструкции казахского текста нейронные сети должны быть обучены таким образом, чтобы распознавать и интерпретировать специфические языковые структуры, культурные оттенки и изменения в языке, произошедшие со временем.

Модели глубокого обучения, используемые для лингвистической реконструкции

В области обработки естественного языка (NLP) глубокие нейронные сети играют ключевую роль в анализе и реконструкции текстов. Для казахского языка, обладающего уникальными особенностями, несколько архитектур, таких как BERT, GPT и T5, продемонстрировали свою эффективность, но каждая из них требует специфической адаптации для достижения наилучших результатов.

Выбор архитектур

BERT (Bidirectional Encoder Representations from Transformers) представляет собой трансформерную модель, которая использует двунаправленный подход к анализу контекста [3]. Это важно для казахского языка, поскольку слова могут иметь разные значения в зависимости от их окружения. BERT эффективен для задач, связанных с выявлением значений и контекстных связей, что делает его особенно подходящим для лингвистической реконструкции. Например, использование BERT для анализа текстов на казахском языке позволяет выделять связи между словами и фразами, что является необходимым для адекватного понимания и интерпретации сложных грамматических конструкций.

GPT (Generative Pre-trained Transformer), в свою очередь, специализируется на генерации текста и предсказании следующих слов на основе контекста. Эта модель может быть применена для задач реконструкции, так как она позволяет восстанавливать недостающие элементы текста и адаптировать его для современного читателя. Применяя GPT к казахским текстам, можно преобразовать классические литературные произведения, сохраняя их смысл, но делая их более доступными для молодежи [4]. Например, если взять произведение Мухтара Ауэзова, GPT может предложить современные синонимы для устаревших слов, улучшая понимание текста.

T5 (Text-to-Text Transfer Transformer) представляет собой универсальный инструмент, который

переводит все задачи обработки текста в формат «текст-в-текст» [5]. Эта модель идеально подходит для разнообразных задач, связанных с лингвистической реконструкцией. В контексте казахского языка T5 может быть использован для преобразования устаревших слов в современные формы или для объяснения сложных понятий, связанных с культурным контекстом. Например, T5 может взять отрывок из классической казахской поэзии и представить его в формате, понятном современному читателю, что значительно расширяет доступность казахской литературы.

### Адаптация моделей

Для эффективной работы с казахским языком необходимо адаптировать модели глубокого обучения. Первый шаг заключается в сборе и подготовке данных. Создание корпуса казахских текстов, включающего разнообразные лексические, морфологические и синтаксические структуры, позволяет моделям более точно обрабатывать и интерпретировать текст. Такой корпус может содержать литературные произведения, архивные документы и современные статьи, обеспечивая всесторонний подход к анализу языка.

Следующий этап — предобучение и дообучение моделей на казахских данных. Этот процесс позволяет моделям лучше адаптироваться к языковым особенностям и культурным контекстам. Например, дообучение модели BERT на корпусе казахских текстов помогает ей глубже понять сложные грамматические структуры и специфические фразеологизмы, характерные для казахского языка. Это приводит к повышению качества интерпретации и реконструкции текстов.

Использование собственных токенизаторов также играет важную роль в адаптации. Казахский язык обладает агглютинативной природой, поэтому стандартные токенизаторы могут не справляться с эффективным разбиением слов на морфемы. Разработка специализированных токенизаторов, учитывающих особенности казахской грамматики, позволяет моделям более точно обрабатывать текст и избегать ошибок, связанных с неправильным пониманием структуры слов.

### Примеры применения глубоких нейронных сетей для казахского текста

Глубокие нейронные сети открывают новые возможности для анализа и реконструкции казахских текстов. Одним из примеров является историко-культурная реконструкция. В текстах, написанных в прошлом, нередко встречаются устаревшие слова и выражения, которые могут быть трудны для современного читателя. Применяя такие модели, как GPT и T5, исследователи смогли перевести эти тексты на современный казахский язык, сохраняя оригинальные культурные и исторические контексты [6]. Например, отрывок из классического произведения может быть преобразован с помощью нейросети, которая добавляет пояснения к устаревшим терминам или заменяет их на более распространенные синонимы. Это значительно облегчает доступность и понимание для широкой аудитории, особенно для молодежи.

Современные тексты также могут быть адаптированы с использованием глубоких нейронных сетей. Например, официальные документы и статьи могут быть переработаны с целью упрощения их восприятия, что особенно важно в образовательной сфере. В таком контексте глубокие нейронные сети могут создавать тексты с понятными формулировками, делая их доступными для разных уровней понимания. Это способствует более широкому распространению информации и помогает учащимся лучше осваивать учебные материалы.

Когнитивная интерпретация и выявление культурных смыслов также являются важными аспектами применения глубоких нейронных сетей. Например, при анализе поэтических текстов нейросети могут выделять символические образы и предлагать интерпретации, что помогает читателям лучше понять культурные контексты произведений. Такой подход не только обогащает читательский опыт, но и способствует более глубокому пониманию казахской культуры.

### Практические результаты и тестирование моделей

Важным аспектом применения глубоких нейронных сетей является оценка их эффективности в контексте казахского языка. В ходе исследований проводились эксперименты, направленные на определение качества реконструированных текстов и точности интерпретации значений. Используемые метрики, такие как точность и полнота, позволяют продемонстрировать, насколько эффективно нейросети справляются с задачами лингвистической реконструкции.

Например, в одном из экспериментов, посвященных реконструкции текстов, связанных с казахскими традициями и обычаями, модели продемонстрировали высокие результаты. В процессе анализа исторических текстов, таких как сказания о Дастане, нейросеть успешно восстанавливала недостающие элементы повествования и предлагала их в новом, более доступном формате. Читатели отмечали, что такие реконструкции не только сохраняли оригинальный смысл, но и обогащали текст современными интерпретациями, что делало его актуальным для новых поколений.

Применение глубоких нейронных сетей также оказалось полезным для создания образовательных материалов. Например, разработанные на базе моделей тексты были полезны для учителей и студентов, позволяя легко находить примеры казахской литературы, адаптированные к образовательным стандартам. Такой подход способствует популяризации казахского языка и литературы, а также помогает молодежи лучше понимать культурное наследие своей страны.

Несмотря на успехи, исследование также выявило некоторые ошибки и ограничения моделей. Например, глубокие нейронные сети могут испытывать трудности в понимании тонкостей казахского языка, что иногда приводит к неправильной интерпретации многозначных слов. Тем не менее, с каждым новым циклом дообучения модели показывают все лучшие результаты, что свидетельствует о том, что, несмотря на существующие сложности, технологии продолжают развиваться и улучшаться.

#### Заключение и перспективы

Подводя итог, можно сказать, что применение глубоких нейронных сетей для лингвистической реконструкции казахских текстов открывает новые возможности для сохранения и популяризации культурного наследия. Современные модели обработки естественного языка успешно адаптируются к уникальным особенностям казахского языка, что позволяет не только анализировать, но и реконструировать тексты с учетом их культурных и исторических контекстов.

Исследования в данной области подчеркивают значимость интеграции технологий глубокого обучения для анализа и интерпретации казахской литературы. Перспективы дальнейшего развития этой области исследований включают улучшение существующих моделей, разработку новых методов для анализа казахского языка и создание интерактивных платформ, которые позволят пользователям глубже взаимодействовать с текстами и способствовать сохранению и развитию казахского языка в цифровую эпоху.

Такое использование нейронных сетей в исследовании казахских текстов не только продвигает развитие технологий, но и укрепляет связь между современными научными достижениями и культурным наследием, делая важный шаг к созданию более просвещенного и информированного общества.

#### Список литературы:

1. Ba L. & Caurana R. Do Deep Nets Really Need to be Deep?, arXiv preprint arXiv:1312.6184, 2013, 521(7553), pp. 1-6.
2. Kalchbrenner N., Grefenstette E., & Blunsom P. A Convolutional Neural Network for Modelling Sentences, In Proc. of the 52nd Annual Meeting of the Association for Computational Linguistics, ACL'2014, Baltimore, MD, USA, 2014, Vol. 1, pp. 655-665.

3. Collobert R., Weston J., Bottou L., Karlen M., Kavukcuglu K., Kuksa P. Natural Language Processing from Scratch, *Journal of Machine Learning Research*, 2011, 12:2493-2537.
4. Santos C.N., Dos & Guimarães V. Boosting Named Entity Recognition with Neural Character Embeddings, *ACL'2015*, pp. 25-33.
5. Malinowski M., Rohrbach M. & Fritz M. 2015. Ask Your Neurons: A Neural based Approach to Answering Questions about Images, *IEEE Int. Conf. on Computer Vision*, 2015, pp. 1-9.
6. Yu L., Hermann K.M., Blunsom P. & Pulman S. Deep Learning for Answer Sentence Selection, *NIPS Deep Learning Workshop*, 2014, 9.