



## ПОЛУАВТОМАТИЧЕСКОЕ ИЗВЛЕЧЕНИЕ ИНФОРМАЦИИ О НАЛИЧИИ ПАРЕЗОВ У ПАЦИЕНТОВ В НЕЙРОХИРУРГИИ ИЗ ИСТОРИЙ БОЛЕЗНИ: ИССЛЕДОВАНИЯ С ПОМОЩЬЮ ТЕХНОЛОГИЙ ИСКУССТВЕННОГО ИНТЕЛЛЕКТА

Г. В. Данилов, А. А. Потапов, М. А. Шифрин, Ю. В. Струнина, К. В. Котик,  
Т. В. Цуканова, Т. Е. Пронкина, Т. А. Ишанкулов, Е. С. Макашова,  
А. В. Косырькова, С. А. Мельченко, Т. Р. Загидуллин

Федеральное государственное автономное учреждение «Национальный медицинский исследовательский центр нейрохирургии имени академика Н. Н. Бурденко» Министерства здравоохранения Российской Федерации, Москва

**РЕЗЮМЕ.** Выявление нежелательных явлений по данным клинических документов необходимо в рамках ретроспективных клинических исследований и в задачах мониторинга безопасности и экономической эффективности медицинской помощи. Поскольку нежелательные явления, как правило, описываются в медицинских записях в виде свободного текста, для извлечения информации о них из десятков тысяч историй болезни требуются специальные технологии.

**ЦЕЛЬ ИССЛЕДОВАНИЯ:** оценить качество предложенного нами алгоритма для полуавтоматической идентификации парезов у пациентов с глиальными опухолями на этапе поступления в нейрохирургический стационар.

**МАТЕРИАЛЫ И МЕТОДЫ:** Для решения данной задачи был применен разработанный нами комплекс технологий с использованием методов искусственного интеллекта (патент RU 2751993С 1). В основе метода лежит отбор специфического для конкретного нежелательного явления лексикона и расшифровка его использования в микроконтекстах. Проанализированы текстовые медицинские записи, первично внесенные врачами с помощью клавиатуры в электронную медицинскую карту «E-med» ФГАУ «НМИЦ нейрохирургии им. ак. Н. Н. Бурденко» Минздрава России в период с 2000 г. по 2017 г.

**РЕЗУЛЬТАТЫ.** Предложенный нами алгоритм позволил выявить парезы на дооперационном этапе с высоким качеством (чувствительность = 0,947, специфичность = 0,965, точность = 0,961, F1-мера = 0,926, ROC AUC = 0,956 [0,941; 0,969]). За счет оптимизации предложенного нами способа (использования только существительных при скрининге словаря) удалось сократить время на его реализацию с 13 до 6 часов практически без потери качества.

**ЗАКЛЮЧЕНИЕ.** Методы анализа текстов, написанных на естественном языке, позволяют улучшить качество извлечения информации из медицинских текстов, что, в частности, может быть успешно применено в исследованиях безопасности оказания нейрохирургической помощи. Проект поддержан грантом РФФИ 18-29-22085. Извлечение данных для глиальных опухолей также было поддержано грантом РФФИ 18-29-01052.

**КЛЮЧЕВЫЕ СЛОВА:** нейрохирургия, тромбоэмболия легочной артерии, осложнения, искусственный интеллект, анализ текстов на естественном языке.

*Для цитирования:* Данилов Г. В., Потапов А. А., Шифрин М. А., Струнина Ю. В., Котик К. В., Цуканова Т. В., Пронкина Т. Е., Ишанкулов Т. А., Макашова Е. С., Косырькова А. В., Мельченко С. А., Загидуллин Т. Р. Полуавтоматическое извлечение информации о наличии парезов у пациентов в нейрохирургии из историй болезни: исследования с помощью технологий искусственного интеллекта. Российский нейрохирургический журнал им. проф. А. Л. Поленова. 2022;14(1–1):52–55

### SEMI-AUTOMATIC INFORMATION EXTRACTION ON THE PRESENCE OF PARESIS IN NEUROSURGICAL PATIENTS FROM HEALTH RECORDS: A RESEARCH USING ARTIFICIAL INTELLIGENCE

G. V. Danilov, A. A. Potapov, M. A. Shifrin, U. V. Strunina, K. V. Kotik,  
T. V. Tsukanova, T. E. Pronkina, T. A. Ishankulov, E. S. Makashova,  
A. V. Kosyrkova, S. A. Melchenko, T. R. Zagidullin

Federal State Autonomous Institution “National Medical Research Center for Neurosurgery named after Academician N. N. Burdenko” of the Ministry of Health of the Russian Federation, Moscow

**SUMMARY.** Adverse events identification in clinical documents is necessary for retrospective clinical research and evaluating medical care safety and cost-effectiveness. Since adverse events are usually reported with free text in medical records, special technologies are required to extract information about them from thousands of medical records.

**MATERIALS AND METHODS:** We introduced a technology to extract information based on natural language processing (patent RU 2751993C1). The method is grounded on preselecting a lexicon specific to a particular adverse event and deciphering its use in microcontexts. The textual medical records, initially typed in by doctors using the keyboard into the electronic health records system «E-med» of the FSAI «NMRC named after ac. N. N. Burdenko» of the Ministry of Health of Russia from 2000 to 2017 were analysed.

**RESULTS.** The technology we proposed enabled us to solve the task of paresis identification with high quality (sensitivity = 0,947, specificity = 0,965, accuracy = 0,961, F1-score = 0,926, ROC AUC = 0,956 [0,941; 0,969]). Optimization of the proposed method (using only nouns when screening a vocabulary) enabled to reduce the time for its implementation from 13 to 6 hours with no major loss in quality.

**CONCLUSION.** Natural language processing can improve the quality of information extraction from medical texts, which, in particular, can be successfully applied in neurosurgical safety research. The project was supported by the RFBR grant 18-29-22085. Data extraction for glial tumors was also supported by RFBR grant 18-29-01052.

**KEY WORDS:** neurosurgery, paresis, complications, artificial intelligence, natural language processing.

*For citation:* Danilov G. V., Potapov A. A., Shifrin M. A., Strunina U. V., Kotik K. V., Tsukanova T. V., Pronkina T. E., Ishankulov T. A., Makashova E. S., Kosyrkova A. V., Melchenko S. A., Zagidullin T. R. Semi-automatic information extraction on the presence of paresis in neurosurgical patients from health records: a research using artificial intelligence. *The Russian Neurosurgical Journal named after prof. A. L. Polenov.* 2022;14(1-1):52-55

**Введение.** Выявление нежелательных явлений по клиническим документам востребовано в ретроспективных научных исследованиях и в задачах мониторинга безопасности и экономической эффективности медицинской помощи [1]. Однако нежелательные явления, как правило, описываются в медицинских записях в виде свободного текста. Поэтому для извлечения информации об осложнениях после нейрохирургических вмешательств, требуется скрупулезная работа экспертов. Извлечение информации из десятков тысяч историй болезни экспертами может оказаться неосуществимым за приемлемый период времени. В качестве средства «усиления» процесса извлечения информации из текстов сегодня можно рассматривать технологии анализа естественного языка (*англ. natural language processing*) и методы искусственного интеллекта.

Наша исследовательская группа предложила алгоритм извлечения информации (АИИ) для обнаружения неблагоприятных событий в нейрохирургии с использованием документов, написанных на естественном морфологически богатом естественном языке [2,3].

**Цель.** Оценить качество работы АИИ для полуавтоматического обнаружения предоперационных парезов у пациентов с глиальными опухолями и его сравнение алгоритмами машинного обучения.

**Материалы и методы:** Набор данных для нашего исследования изначально был подготовлен в рамках проекта по прогнозированию мышечной силы методом глубокого обучения с использованием медицинских изображений (при поддержке гранта РФФИ 19-29-01154). Техническая задача настоящей работы заключалась в классификации пациентов как имеющих слабость конечностей (парез) или не имеющих ее на дооперационном этапе по данным описания неврологического статуса в форме произвольного неструктурированного текста. Предоперационные текстовые данные и разделы выписных эпикризов, относящиеся к предоперационному периоду, были пер-

вично набраны врачами с помощью клавиатуры и извлечены из электронной медицинской карты «E-med» ФГАУ «НМИЦ нейрохирургии им. ак. Н. Н. Бурденко» Минздрава России для когорты из 1167 больных (средний возраст  $39,7 \pm 17,8$  лет, 578 (49,5 %) мужчин), перенесших резекцию глиальных опухолей в НМИЦ нейрохирургии им. Н. Н. Бурденко в период с 2009 года по 2018 год. Массив данных оценивали и аннотировали эксперты с использованием специально разработанного для этого приложения. Каждый из двух независимых экспертов-нейрохирургов определял наличие или отсутствие пареза, третий независимый эксперт на следующем шаге устранял разногласия между оценками экспертов. В результате формировали целевую бинарную переменную, принимавшую значения 0 (парез в дооперационном периоде отсутствует) и 1 (в предоперационном периоде выявлен парез). Размеченный набор данных использовался для тестирования качества АИИ в небольшой модификации. Все данные, извлеченные из электронной медицинской карты, были предобработаны и проанализированы с помощью языка программирования R (версия 3.5.0) в среде RStudio Server IDE (версия 1.2.519) с использованием библиотек *rsample*, *drlib*, *broom*, *yardstick*, *glmnet*, *doMC*, *dplyr*, *накеты*, *tidyr*, *stringr*, *tidytext*, *tidyverse*, *quanteda*, *qdapRegex*, *tm*, *scales* и *cutpointr*.

Исходные текстовые документы, содержащие описание неврологического статуса пациентов при поступлении, были предобработаны следующим образом: тексты разделены на слова по пробелу, слова переведены в нижний регистр, удалены все символы, кроме буквенно-цифровых символов и одиночных пробелов, удалены стоп-слова и бессмысленные слова (одиночные буквы, артефакты и т.д.), исправлены орфографические ошибки, слова приведены к начальной форме. Далее эти слова просматривал эксперт и отбирал те, что с высокой вероятностью могут быть использованы в описании пареза. Контекст каждого такого отобранного слова далее определялся

по окружающим его терминам в разных фразах. Таким образом, для каждой из фраз, содержащей какое-либо отобранное слово, определяли оценку в 1 балл (при явном указании фразы на парез конечности), 0 баллов (если о парезе не шла речь) и 0,5 балла (при невозможности сделать заключение о парезе с определенностью). По совокупности фраз, оцененных на 1, 0 или 0,5 баллов в тексте каждой истории болезни делали заключение о наличии или отсутствии пареза до операции для соответствующего клинического случая. Наличие пареза устанавливалось для тех историй болезни, в которых хотя бы одна фраза была оценена на 1 балл.

#### Результаты.

Мы наблюдали высокий уровень согласия между двумя экспертами (коэффициент каппа Коэна = 0,852), аннотировавшими набор данных. Тем не менее потребовалась коррекция разметки третьим экспертом. В результате скрининга общего словаря из 9932 слов, извлеченных из текстов, были отобраны 738 терминов, потенциально имевших отношение к описанию парезов. Далее эти слова были найдены в 7413 фразах из 5 слов, которые были размечены экспертом на группы в 1, 0 или 0,5 баллов. На весь процесс работы с текстами потребовалось около 13 часов.

Предложенный нами алгоритм позволил выявить парезы на дооперационном этапе с высоким качеством (чувствительность = 0,947, специфичность = 0,965, точность = 0,961, F1-мера = 0,926, ROC AUC = 0,956 [0,941; 0,969]). За счет оптимизации предложенного нами способа (использования только существенных при скрининге словаря) удалось сократить время на его реализацию до 6 часов практически без потери качества (чувствительность = 0,944, специфичность = 0,965, точность = 0,960, F1-мера = 0,924, ROC AUC = 0,954 [0,938; 0,969]) [4,5]. Более подробные сведения о полученных результатах исследований представлены в устном сообщении на конференции «Поленовские чтения» 2022 г.

**Заключение.** Исследование неблагоприятных событий в медицине, в том числе — с использованием методов анализа текстов, способствует повышению безопасности оказания медицинской помощи [6]. Извлечение информации из клинических текстов с помощью технологий искусственного интеллекта ранее применяли в различных областях медицины (например, при анализе инфекционной безопасности или при испытаниях лекарств) [3]. Однако эти технологии редко использовали для выявления нежелательных явлений в нейрохирургии [7–9]. Мы предложили и протестировали алгоритм (отличный от методов, основанных на математических моделях), который может быть применим для разметки клинических текстов и обнаружения нежелательных явлений в качестве первого шага к последующей автоматизации этого процесса с помощью машинного обучения. Нам удалось снизить нагрузку на эксперта с помощью нашего алгоритма примерно на 50 %, сохранив качество обнаружения парезов на стабильно высоком уровне.

Таким образом, методы анализа текстов, написанных на естественном языке, позволяют улучшить качество извлечения информации из медицинских текстов, что, в частности, может быть успешно применено в исследованиях безопасности оказания нейрохирургической помощи.

**Ограничения исследований.** Ограничения наших исследований связаны с лимитами информации, представленной в медицинской документации и субъективностью оценок врачей в описании явлений произвольным текстом.

**Конфликт интересов.** Авторы заявляют об отсутствии конфликта интересов. **Conflict of interest.** The author declares no conflict of interest.

**Финансирование.** Работа поддержана Российским фондом фундаментальных исследований (грант 18–29–22085).

**Financing.** The study was supported by the Russian Foundation for Basic Research (grant 18–29–22085).

**Соблюдение прав пациентов и правил биоэтики:** Все пациенты подписали информированное согласие на участие в исследовании. **Compliance with patient rights and principles of bioethics.** All patients gave written informed consent to participate in the study

#### ORCID авторов / ORCID of authors:

Данилов Глеб Валерьевич/Danilov Gleb Valerievich  
<https://orcid.org/0000-0003-1442-5993>

Потапов Александр Александрович/  
Potapov Alexander Aleksandrovich  
<https://orcid.org/0000-0001-8343-3511>

Шифрин Михаил Абрамович/Shifrin Mikhail Abramovich  
<https://orcid.org/0000-0002-9606-5559>

Струнина Юлия Владимировна/  
Strunina Yuliya Vladimirovna  
<https://orcid.org/0000-0001-5010-6661>

Котик Константин Владимирович/  
Kotik Konstantin Vladimirovich  
<https://orcid.org/0000-0002-3524-4983>

Цуканова Татьяна Васильевна/  
Tsukanova Tatyana Vasilievna  
<https://orcid.org/0000-0002-0046-1312>

Пронкина Татьяна Евгеньевна/  
Pronkina Tatyana Yevgenievna  
<https://orcid.org/0000-0003-0068-4599>

Ишанкулов Тимур Александрович/  
Ishankulov Timur Aleksandrovich  
<https://orcid.org/0000-0002-6509-4242>

Макашова Елизавета Сергеевна/  
Makashova Yelizaveta Sergeevna  
<https://orcid.org/0000-0003-2441-8818>

Косырькова Александра Вячеславовна/  
Kosyrkova Aleksandra Vyacheslavovna  
<https://orcid.org/0000-0002-3019-5203>

Мельченко Семен Андреевич/  
Melchenko Semen Andreyevich  
<https://orcid.org/0000-0001-7060-0667>

Загидуллин Тимур Рустамович/  
Zagidullin Timur Rustamovich  
<https://orcid.org/0000-0002-5621-6834>

## Литература/References

1. Wang Y, Wang L, Rastegar-Mojarad M, et al. Clinical information extraction applications: A literature review. *J Biomed Inform.* 2018;77:34–49. doi:10.1016/j.jbi.2017.11.011
2. Danilov G, Shifrin M, Strunina U, Pronkina T, Potapov A. An Information Extraction Algorithm for Detecting Adverse Events in Neurosurgery Using Documents Written in a Natural Rich-in-Morphology Language. *Stud Health Technol Inform.* 2019;262:194–197. doi:10.3233/SHTI190051
3. Данилов Г.В., Шифрин М.А., Потапов А.А., et al. Способ извлечения информации из неструктурированных текстов, написанных на естественном языке (патент RU 2751993 C 1). Опубликовано в сети 21 июля 2021 г. По состоянию на 27 декабря 2021 г. [https://www.fips.ru/registers-doc-view/fips\\_servlet?DB=RUPAT&DocNumber=2751993&TypeFile=html](https://www.fips.ru/registers-doc-view/fips_servlet?DB=RUPAT&DocNumber=2751993&TypeFile=html) [Danilov G. V., Shifrin M. A., Potapov A. A., et al. Method for extracting information from unstructured texts written in natural language (patent RU 2751993 C 1). Published online July 21, 2021. Accessed December 27, 2021. [https://www.fips.ru/registers-doc-view/fips\\_servlet?DB=RUPAT&DocNumber=2751993&TypeFile=html](https://www.fips.ru/registers-doc-view/fips_servlet?DB=RUPAT&DocNumber=2751993&TypeFile=html) (In Russ.).]
4. Danilov G, Shifrin M, Strunina Y, et al. Detection of muscle weakness in medical texts using natural language processing. In: *Studies in Health Technology and Informatics*. Vol 270. IOS Press; 2020:163–167. doi:10.3233/SHTI200143
5. Danilov G, Shifrin M, Strunina Y, et al. Semiautomated approach for muscle weakness detection in clinical texts. In: *Studies in Health Technology and Informatics*. Vol 272. IOS Press; 2020:55–58. doi:10.3233/SHTI200492
6. Young IJB, Luz S, Lone N. A systematic review of natural language processing for classification tasks in the field of incident reporting and adverse event analysis. *Int J Med Inform.* 2019;132:103971. doi:10.1016/j.ijmedinf.2019.103971
7. Gaebel J, Kolter T, Arlt F, Denecke K. Extraction Of Adverse Events From Clinical Documents To Support Decision Making Using Semantic Preprocessing. *Stud Health Technol Inform.* 2015;216:1030.
8. Campillo-Gimenez B, Garcelon N, Jarno P, Chaplain JM, Cuggia M. Full-text automated detection of surgical site infections secondary to neurosurgery in Rennes, France. *Stud Health Technol Inform.* 2013;192:572–575.
9. Tvardik N, Kergourlay I, Bittar A, Segond F, Darmoni S, Metzger M-H. Accuracy of using natural language processing methods for identifying healthcare-associated infections. *Int J Med Inform.* 2018;117:96–102. doi:10.1016/j.ijmedinf.2018.06.002