ISSN 2070-4909 (print) ISSN 2070-4933 (online)

Фармакоэконо Современная фармакоэкономика и фармакоэпидемиология



FARMAKOEKONOMIKA

Modern Pharmacoeconomics and Pharmacoepidemiology

2025 Vol. 18 No. 3

нтах можно получить в редакции. Тел.: +7 (495) 649-54-95; эл. почта: info@irbis-1.ru Том **₹8**

Сверточные нейронные сети и визуальные трансформеры в диагностике опухолей кожи: сравнительный анализ эффективности моделей искусственного интеллекта в программах компьютерного зрения

А.И. Ламоткин 1,2 , Д.И. Корабельников 1

Для контактов: Андрей Игоревич Ламоткин, e-mail: lamotkin.an@mail.ru

РЕЗЮМЕ

Актуальность. При применении программ компьютерного зрения в моделях искусственного интеллекта (ИИ) для целей диагностики кожных опухолей, особенно меланомы, критически важны минимальные ошибки классификации. Новая архитектура моделей – визуальные трансформеры (англ. vision transformer, ViT) демонстрируют перспективные результаты, однако их эффективность для классификации кожных новообразований изучена недостаточно. Данное исследование является одним из первых, где проведено прямое сравнение сверточных нейронных сетей (англ. convolutional neural network, CNN) и ViT на клинически значимых метриках, что особенно важно для внедрения ИИ в дерматологическую и онкологическую практику.

Цель: сравнить эффективность CNN и ViT в задачах бинарной классификации опухолей кожи: «меланома / не меланома» и «доброкачественное/злокачественное».

Материал и методы. В исследовании сравнивалась эффективность программ компьютерного зрения (приложения для смартфона) с использованием архитектуры CNN (MobileNetV2, Xception) и трансформера ViT. Тестирование проводилось на независимых наборах данных (3000 и 4800 изображений соответственно) с оценкой по метрикам точности, чувствительности, специфичности. Применялись методы аугментации. балансировки классов, оптимизации гиперпараметров и очистки данных от артефактов.

Результаты. ViT показал превосходство над CNN: в задаче «меланома / не меланома» точность составила 92,93% против 88% у Xception, а в задаче «доброкачественное/злокачественное» - 91,35% против 85%. Модель трансформера продемонстрировала лучшую специфичность (до 95%).

Заключение. ViT обеспечивает более высокую точность за счет анализа глобальных паттернов, но требует тщательной настройки и качественных данных. CNN остаются стабильным решением при ограниченных данных. Для клинического применения рекомендовано комбинирование обеих архитектур, позволяющее повысить эффективность диагностики.

КЛЮЧЕВЫЕ СЛОВА

искусственный интеллект, компьютерное зрение, сверточные нейронные сети, трансформеры, классификация изображений, меланома, доброкачественные опухоли, злокачественные опухоли, диагностика

Для цитирования

Ламоткин А.И., Корабельников Д.И. Сверточные нейронные сети и визуальные трансформеры в диагностике опухолей кожи: сравнительный анализ эффективности моделей искусственного интеллекта в программах компьютерного зрения. ФАРМА-КОЗКОНОМИКА. Современная фармакоэкономика и фармакоэпидемиология. 2025; 18 (3): 365-375. https://doi.org/10.17749/ 2070-4909/farmakoekonomika.2025.327.

 $^{^{1}}$ Автономная некоммерческая организация дополнительного профессионального образования «Московский медико-социальный институт им. Ф.П. Гааза» (2-я Брестская ул., д. 5, Москва 123056, Российская Федерация)

² Федеральное государственное бюджетное учреждение «Центральный научно-исследовательский институт организации и информатизации здравоохранения» Министерства здравоохранения Российской Федерации (ул. Добролюбова, д. 11, Москва 127254, Российская Федерация)



Convolutional neural networks and transformers in skin tumor diagnostics: a comparative analysis of the efficiency of artificial intelligence models in computer vision programs

A.I. Lamotkin^{1,2}, D.I. Korabelnikov¹

- ¹ Moscow Haass Medical and Social Institute (5 2nd Brestskaya Str., Moscow 123056, Russian Federation)
- ² Central Research Institute of Organization and Informatization of Healthcare (11 Dobrolyubov Str., Moscow 127254, Russian Federation)

Corresponding author: Andrey I. Lamotkin, e-mail: lamotkin.an@mail.ru

ABSTRACT

Background. When applying computer vision programs in artificial intelligence (AI) models to diagnose skin tumors, melanoma in particular, even minimal classification errors are critical. Vision transformers (ViT), a new model architecture, have shown promising results in computer vision; however, their efficiency in classifying skin lesions has received insufficient research attention. This study is one of the first to directly compare convolutional neural networks (CNNs) and ViT on clinically relevant metrics, which is especially important for the implementation of AI in dermatological and oncological practice.

Objective: To compare the performance of CNNs and ViT in the tasks of binary classification of skin lesions: "melanoma/not melanoma" and "benign/malignant".

Material and methods. The study used CNN (MobileNetV2, Xception) and ViT architectures. Testing was carried out on independent datasets (3000 and 4800 images, respectively) with assessment by the metrics of accuracy, sensitivity, specificity. Augmentation, class balancing, hyperparameter optimization, and data cleaning from artifacts were used.

Results. ViT showed superiority over CNNs. Thus, in the "melanoma/non-melanoma" task, the accuracy was 92.93% versus 88% for Xception, and in the "benign/malignant" task – 91.35% versus 85%. Transformer model demonstrated better specificity (up to 95%).

Conclusion. ViT provides higher accuracy due to the analysis of global patterns, although requiring careful tuning and high-quality data. CNNs remain a stable solution with limited data. For clinical use, a combination of both architectures is recommended to improve the reliability of diagnostics.

KEYWORDS

artificial intelligence, computer vision, convolutional neural networks, visual transformers, image classification, melanoma, benign neoplasms, malignant neoplasms, diagnosis

For citation

Lamotkin A.I., Korabelnikov D.I. Convolutional neural networks and transformers in skin tumor diagnostics: a comparative analysis of the efficiency of artificial intelligence models in computer vision programs. *FARMAKOEKONOMIKA. Sovremennaya farmakoekonomika i farmakoepidemiologiya / FARMAKOEKONOMIKA. Modern Pharmacoeconomics and Pharmacoepidemiology.* 2025; 18 (3): 365–375 (in Russ.). https://doi.org/10.17749/2070-4909/farmakoekonomika.2025.327.

Основные моменты

Что уже известно об этой теме?

- Сверточные нейронные сети (англ. convolutional neural networks, CNNs) широко применяются в медицинской диагностике, но имеют ограничения в интерпретируемости и устойчивости к сдвигам домена
- Визуальные трансформеры (англ. vision transformer, ViT) демонстрируют высокую эффективность в задачах классификации изображений
- Балансировка классов и аугментация данных критически важны для качества моделей

Что нового дает статья?

- Данное исследование является одним из первых в мире с прямым сравнением архитектур CNN и ViT в задаче классификации кожных опухолей с использованием клинически значимых метрик
- ViT показал превосходство в точности (до 92,93%) и специфичности (до 95%), особенно в сложных диагностических случаях

Как это может повлиять на клиническую практику в обозримом будущем?

- ▶ Внедрение ViT повысит точность диагностики меланомы и снизить количество ложноположительных результатов
- Гибридные решения (CNN + ViT) могут стать стандартом для автоматизированной дифференциальной диагностики доброкачественных и злокачественных опухолей кожи, особенно меланомы
- Результаты подчеркивают важность качества данных и методов предобработки для внедрения искусственного интеллекта в медицину

Highlights

What is already known about the subject?

- Convolutional neural networks (CNNs) are widely used in medical diagnostics; however, these architectures have limitations in interpretability and robustness to domain shifts
- Vision transformers (ViT) demonstrate high performance in image classification tasks
- ▶ Class balancing and data augmentation are critical for model quality

What are the new findings?

- ▶ This study is one of the first to directly compare CNN and ViT architectures in the task of skin tumor classification using clinically relevant metrics
- ViT showed superiority in accuracy (up to 92.93%) and specificity (up to 95%), particularly in complex diagnostic cases

How might it impact the clinical practice in the foreseeable future?

- ➤ ViT implementation will improve the accuracy of melanoma diagnostics and reduce the number of false positive results
- Hybrid solutions (CNN + ViT) can become the standard for automated differential diagnosis of benign and malignant skin lesions, melanoma in particular
- ► The results highlight the importance of data quality and pre-processing methods for implementing artificial intelligence in medicine

ВВЕДЕНИЕ / INTRODUCTION

Компьютерное зрение (англ. computer vision, CV) — это область искусственного интеллекта (ИИ), направленная на автоматизированную обработку и анализ визуальных данных, включая 2D- и 3D-изображения, а также видеопотоки [1–5]. Технологии CV находят применение в разнообразных сферах, таких как нефтегазовая промышленность, агропромышленный комплекс [1–5], медицинская диагностика [6, 7] и многие другие.

Основные задачи CV традиционно подразделяются на три ключевые категории:

- классификация изображений отнесение изображения к определенному классу на основе его содержания [8];
- детекция изображений идентификация и локализация объектов на изображении [8];
- сегментация изображений разделение изображения на семантически значимые области [8].

Классификация изображений занимает особое место благодаря своей актуальности и востребованности, особенно в медицинской диагностике [9]. Данная задача обычно формулируется в рамках обучения с учителем (англ. supervised machine learning), где модель обучается предсказывать метку y (целевой класс) на основе набора признаков X, извлеченных из изображения. Математически цель классификации заключается в нахождении функции h(x), которая устанавливает соответствие между входным изображением X и его меткой y.

Сверточные нейронные сети (англ. convolutional neural network, CNN) — это разновидность алгоритма глубокого обучения и ключевая технология, лежащая в основе современной области ИИ [10]. Их архитектура включает сверточные слои с нелинейными функциями активации, слои подвыборки (пулинга) и полносвязные слои, что позволяет эффективно выявлять сложные паттерны в визуальных данных [11]. Благодаря этим свойствам СNN нашли широкое применение в обработке

медицинских изображений, повышая точность и скорость диагностики.

СNN активно используются для решения таких задач, как классификация, сегментация и детекция изображений, полученных с применением различных диагностических методов (эндоскопия, рентгенография, магнитно-резонансная и компьютерная томография, ультразвуковое исследование), а также в дерматологии и гистопатологии¹ [12–15]. На рисунке 1 показан пример использования CNN для анализа изображений опухолей кожи. С помощью этих алгоритмов успешно диагностируются переломы костей скелета, пневмонии, онкологические заболевания, а также прогнозируется течение рака и выполняется классификация мутаций на основе генетических данных [16–21].

Однако, несмотря на высокую эффективность, CNN обладают рядом ограничений. Одна из ключевых проблем — недостаточная интерпретируемость: эти модели часто функционируют как «черные ящики», не предоставляя понятного объяснения своих заключений. Для ее преодоления разрабатываются методы визуализации, такие как градиентное картирование активации классов (англ. gradient-weighted class activation mapping, Grad-CAM) — метод в глубоком обучении, визуализирующий области изображения, на которых фокусируется модель при прогнозах, что помогает понять, как она принимает решения. Но проблема остается актуальной, особенно в медицине, где обоснованность диагноза критически важна [23].

Еще одним существенным ограничением является чувствительность CNN к сдвигам домена: их точность может значительно снижаться при обработке данных, отличающихся от обучающей выборки, — например, изображений, полученных на другом оборудовании или в ином медицинском учреждении [24—26]. Это ставит под вопрос надежность и универсальность таких систем в реальной клинической практике.

Таким образом, несмотря на значительный потенциал CNN в медицинской диагностике, их внедрение требует решения

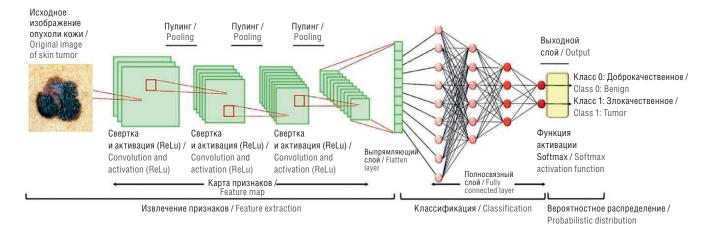


Рисунок 1. Использование сверточных нейронных сетей для анализа дерматологических изображений (адаптировано из [22]).
ReLU (англ. rectified linear unit) — нелинейная функция активации в нейронных сетях, которая широко используется в глубоком обучении и преобразует входное значение в значение от 0 до положительной бесконечности

 $\textbf{Figure 1.} \ Convolutional \ neural \ networks \ in \ dermatological \ image \ analysis \ (adapted \ from \ [22]).$

ReLU – rectified linear unit (a nonlinear activation function in neural networks that is widely used in deep learning, transforming the input value into a value between 0 and positive infinity)

¹ Ламоткин А.И., Ламоткин И.А., Корабельников Д.И. Программа для визуальной идентификации злокачественных и доброкачественных опухолей кожи «Melanoma check». Свидетельство о государственной регистрации программы для ЭВМ № RU 2024668565, заявка от 08.08.2024; Ламоткин А.И., Ламоткин И.А., Корабельников Д.И. Программа для визуальной идентификации злокачественных и доброкачественных опухолей кожи «Derma Onko Check». Свидетельство о государственной регистрации программы для ЭВМ № RU 2024668566, заявка от 08.08.2024.

проблем объяснимости и устойчивости к изменениям входных данных.

Визуальные трансформеры (англ. vision transformer, ViT) это класс моделей глубокого обучения, адаптирующих архитектуру визуальных трансформеров, изначально разработанных для обработки естественного языка (англ. natural language processing, NLP), к задачам анализа изображений [27]. Их адаптация к визуальным задачам привела к появлению преобразователей зрения, которые обрабатывают изображения как последовательность фрагментов (патчей, англ. patch), применяя к ним трансформерную архитектуру [28] (рис. 2). Каждый блок трансформера состоит из двух слоев: слой внимания (англ. attention layer) и слой с прямой связью, применяемый вдоль измерения признаков, - простейший односвязный слой (англ. feed-forward layer). В отличие от CNN, ViT исключает индуктивные смещения, связанные с локальными связями, и напрямую изучает пространственные иерархии через взаимодействие патчей, что обеспечивает большую гибкость при работе с разнородными визуальными данными [29, 30]. Высокую эффективность ViT в CV часто объясняют дизайном слоев, называемым «многоголовое внимание» (англ. multi-headed attention).

Данный подход продемонстрировал конкурентоспособные результаты, превзойдя CNN на ключевых эталонных наборах данных [28, 31–33]. Для повышения интерпретируемости ViT разрабатываются методы визуализации, такие как анализ карт внимания, которые выявляют значимые области изображения, используемые моделью для принятия решений [34–35].

Однако внедрение ViT в медицинскую практику сопряжено с рядом ограничений. Во-первых, их эффективное обучение требует значительных объемов размеченных данных, что часто недостижимо в медицине из-за ограниченности выборок [36—38]. Во-вторых, относительная новизна архитектуры означает отсутствие устоявшихся методик оптимизации и валидации, в отличие от хорошо изученных CNN.

Хотя гибридные модели, сочетающие свертки и механизмы внимания, активно исследуются, прямые сравнения эффективности CNN и ViT проводятся редко из-за методологических сложностей (например, различий в подходах к предобработке данных или настройке гиперпараметров) и новизны использования [39–42].

В связи с этим ключевой вопрос для медицинской визуализации формулируется следующим образом: «Какая архитектура (CNN или ViT) обеспечивает оптимальные результаты при работе с медицинскими изображениями, учитывая специфику данных и требования к интерпретируемости?» В настоящей статье представлен сравнительный анализ эффективности моделей на основе CNN и ViT в программах CV для ЭВМ (приложения для смартфонов) в задачах классификации изображений кожи. Данная научная работа является одной из первых в мировой научной литературе, где проводится прямое сравнение CNN и ViT в диагностике опухолей кожи с применением программ CV.

Цель – сравнить эффективность CNN и ViT в задачах бинарной классификации опухолей кожи: «меланома / не меланома» и «доброкачественное/злокачественное».

MATEРИАЛ И METOДЫ / MATERIAL AND METHODS

Объекты исследования / Study objects

В исследовании рассматривались две версии модели CNN, основанные на архитектурах MobileNetV2 и Хсерtion (обе были предобучены на наборе данных ImageNet-1K), а также модель визуального трансформера на основе Google ViT (предобученная на ImageNet21k).

Описания моделей / Model descriptions

Модель CNN 1-й версии

В 1-й версии модели CNN по классификации «меланома / не меланома» применена легкая архитектура MobileNetV2, оптимизированная для мобильных устройств. Она использует разделяемые по глубине (англ. depthwise-separable) свертки (3×3 для пространственной фильтрации и 1×1 для изменения числа каналов) и включает около 53 сверточных слоев. Такая компактность дает быстрый вывод и экономию ресурсов, но при небольшом числе параметров сложно распознавать очень тонкие или сложные паттерны, например атипичную пигментацию или мельчайшие границы меланомы.

В 1-й версии модели CNN по классификации «доброкачественное/злокачественное» использована базовая реализация Хсерtion с датасетом из 20 тыс. изображений (НАМ10000 с аугментацией) и 240 гистологически подтвержденных случаев, взятых из базы данных профессора И.А. Ламоткина.

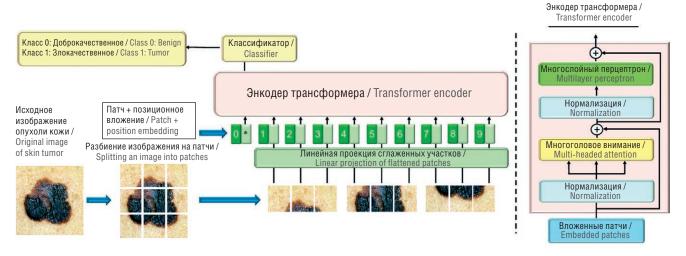


Рисунок 2. Использование визуальных трансформеров для анализа дерматологических изображений (адаптировано из [22])

Figure 2. Use of computer vision transformers for dermatological image analysis (adapted from [22])

Модель CNN 2-й версии

Во 2-й версии модели CNN по классификации «меланома / не меланома» использована архитектура Хсерtion с примерно 71 сверточным слоем, 36 разделяемыми по глубине блоками, включающими входной, средний и выходной потоки (англ. entry, middle и exit flow соответственно)², и внешними пропускными соединениями — обходными путями (англ. skip connections)³.

Во 2-й версии модели CNN по классификации «доброкачественное/злокачественное» сохранена архитектура Xception, но добавлены дополнительные сверточные слои для улучшения детекции текстурных особенностей, что позволило точнее выявлять мелкие детали образований.

В версии Xception модели реализован расширенный конвейер аугментации, включающий:

- коррекцию освещения адаптивную гистограммную эквализацию (англ. contrast limited adaptive histogram equalization, CLAHE) для усиления контраста в условиях затенения или пересвета;
- геометрические трансформации случайные горизонтальные сдвиги (±15% от ширины изображения) и вертикальные отражения для имитации вариаций положения кожи при съемке;
- цветовые искажения добавление шума в RGB-каналы 4 (σ =0,1) и изменение насыщенности (\pm 20%) для повышения устойчивости к артефактам камеры.

За счет более глубокой сепарации пространственных и канальных признаков и увеличенного числа параметров модель Хсерtion стала лучше фиксировать мелкомасштабные детали и асимметрию границ пигментных образований.

В версии Хсерtion модели реализованы ключевые улучшения: – расширенный датасет (увеличение до 45 тыс. изображений (НАМ10000) и 380 верифицированных случаев);

- усовершенствованная предобработка:
 - сложные методы нормализации контраста и яркости,
 - механизмы шумоподавления,
- контролируемые изменения освещения и углов поворота;
- оптимизация обучения:
- скорость обучения (параметр, определяющий величину шага изменения весов, англ. learning rate) 0,0001;
- ранняя остановка при отсутствии увеличения точности на валидации (три эпохи);
- очистка данных (удаление изображений с артефактами (черные края)).

Модель ViT

Модель на основе ViT представляет собой архитектуру глубокого обучения, адаптированную для задач CV. Она разбивает изображение на последовательность патчей фиксированного размера (16×16 пикселей), которые преобразуются в эмбеддинги (англ. embedding)⁵ и обрабатываются трансформерными блоками с механизмом самовнимания (англ. self-attention)⁶.

Для обучения использовалась модель Google ViT с размером патча 16×16 пикселей, предобученная на наборе данных ImageNet21k, что обеспечило высокую обобщающую способность благодаря большому объему разнообразных изображений.

Классификационная голова (англ. head)⁷ модели была заменена и адаптирована для задачи классификации кожных опухолей, что позволило оптимизировать модель для бинарных задач «меланома / не меланома» и «доброкачественное/ злокачественное».

В исследовании применялась конфигурация ViT с частичным замораживанием первых семи блоков, настройкой последних четырех блоков и классификатора, скоростью обучения 1е-4 с косинусным расписанием, размером батча (англ. batch — пакет данных) 64 и аугментацией (горизонтальные отражения, коррекция яркости/контраста). Для повышения интерпретируемости использовались карты внимания (англ. attention maps), которые визуализируют области изображения, наиболее значимые для классификации. Модель обучалась на сбалансированных наборах данных.

Подготовка данных / Data preparation

Для задачи «меланома / не меланома» тренировочная выборка была сформирована с учетом балансировки классов. Оптимальное соотношение достигалось при формировании класса «не меланома» из 50% доброкачественных меланоцитарных опухолей кожи и 50% других доброкачественных опухолей, что обеспечивало максимальную диагностическую точность.

Для задачи бинарной классификации «доброкачественное/ злокачественное» использовалось сбалансированное распределение подклассов внутри каждой категории (например, равные пропорции различных типов как доброкачественных, так и злокачественных опухолей), что способствовало устойчивости модели и минимизации систематических ошибок классификации.

Тестирование / Testing

Две версии модели CNN, основанные на архитектурах MobileNetV2 и Xception, сравнивались между собой по эффективности в классификации опухолей кожи по классам «меланома / не меланома», а также «доброкачественное/зло-качественное».

Проведено сравнительное тестирование двух версий модели CNN на едином тестовом наборе из 3000 изображений (1500 — «меланома», 1500 — «не меланома»), исключенных из обучающих и валидационных данных, но сохраненных при одинаковых параметрах съемки (освещение, ракурс, разрешение).

Выполнено сравнительное тестирование двух версий модели CNN на едином тестовом наборе из 4800 изображений (2400 — «доброкачественное», 2400 — «злокачественное»), исключенных из обучающих и валидационных данных, но

² Структурные компоненты архитектуры Хсерtion: Entry flow отвечает за начальную обработку входных данных, Middle flow – за углубленную обработку признаков, а Exit flow – за финальное извлечение признаков перед классификацией.

³ Skip connections – техника в глубоком обучении, которая позволяет информации «перепрыгивать» через один или несколько слотов в нейронной сети, что помогает избегать проблем исчезающего градиента и сохранять информацию между ними.

⁴ RGB (англ. red, green, blue) – аддитивная цветовая модель, в которой цвета синтезируются путем добавления трех цветов – красного, зеленого, синего в различных количествах.

⁵ Способ представления объектов (слов, предложений, изображений и других типов данных) в виде числовых векторов, которые позволяют нейронным сетям работать с данными и анализировать их взаимосвязи.

⁶ Тип механизма внимания, используемый в моделях машинного обучения, с помощью которого оценивают важность токенов, или слов во входной последовательности, чтобы лучше понимать отношения между ними.

⁷ Компонент, который классифицирует представления объектов по нескольким категориям.

сохраненных при одинаковых параметрах съемки (освещение, ракурс, разрешение).

По тем же критериям проведено сравнение обеих версий моделей CNN с моделью с архитектурой на основе ViT.

Для объективизации и воспроизводимости сравнительного анализа во всех тестированиях использовался единый набор изображений, не входивших в обучающую и валидационную выборки. Все тестовые данные соответствовали одинаковым параметрам съемки (освещение, ракурс, разрешение) и были равномерно распределены между классами. Оценка проводилась с помощью матрицы ошибок, метода Grad-CAM для визуализации областей интереса модели, а также стандартных метрик — точность, чувствительность, специфичность, площадь под кривой рабочей характеристики приемника (англ. area under curve of receiver operating characteristic, AUC ROC) и F1-score.

РЕЗУЛЬТАТЫ / RESULTS

Сравнение двух версий CNN / Comparison of two CNN versions

Классификация изображений по классам «меланома / не меланома»

Обе модели CNN реализовали бинарную классификацию со сбалансированными классами (Softmax, два выходных нейрона), однако архитектурные и методологические различия определили значимую разницу в метриках.

Ключевыми факторами улучшения для модели CNN 2-й версии стали:

- углубленная архитектура (дополнительные сверточные слои улучшили детекцию текстур);
- расширенная аугментация (новые методы трансформации изображений);
- оптимизация параметров (точная настройка скорости обучения:
- качество данных (очистка от артефактов и увеличение выборки);
- контроль обучения (ранняя остановка предотвратила переобучение).

Модель Хсерtion продемонстрировала превосходство над моделью MobileNetV2: точность повысилась с 79% до 88%, специфичность — с 67% до 87%, AUC ROC — с 86% до 95% при сохранении высокой чувствительности (89–90%).

Ключевыми факторами улучшения для модели Xception стали:

- архитектурная оптимизация (глубокая обработка признаков Xception обеспечила лучшее распознавание паттернов меланомы);
- гиперпараметрическая настройка (эксперименты со скоростью обучения выявили его критическое влияние на сходимость модели, тогда как коррекция размера батча (англ. batch size)⁸ и масштабирование изображений минимизировали необходимость переобучения);
- расширенная аугментация (добавление вариаций освещенности, горизонтальных сдвигов и синтетического увеличения датасета повысило инвариантность модели к артефактам);
- ранняя остановка (мониторинг валидационной точности с остановкой обучения при отсутствии улучшений в течение трех эпох предотвратил деградацию модели).

Сравнение матрицы ошибок для двух версий моделей CNN по классификации «меланома / не меланома» представлено на рисунках 3а, 3b.

Классификация изображений по классам «доброкачественное/злокачественное»

Обе версии реализовали бинарную классификацию на базе архитектуры Xception со сбалансированными классами (Softmax, два выходных нейрона), однако методологические улучшения во 2-й версии определили значимый рост метрик.

Модель CNN 1-й версии при тестировании показала точность 81%, чувствительность 84%, специфичность 78%, AUC ROC 90%.

Модель CNN 2-й версии продемонстрировала превосходство над первоначальной реализацией: точность повысилась с 81% до 85%, чувствительность — с 84% до 87%, специфичность — с 78% до 83%, AUC ROC — с 90% до 93%. F1-score для злокачественных образований вырос с 0.82 до 0.86.

Сравнение матриц ошибок двух версий моделей CNN по классификации «доброкачественное/злокачественное» представлено на рисунках 4а, 4b.

Сравнение моделей CNN и ViT / Comparison of CNN and ViT models

ViT продемонстрировал преимущество в выявлении глобальных паттернов (асимметрия, распределение пигмента), но требовал на 15–20% больше данных по сравнению с CNN. Ранняя остановка после двух эпох без улучшения валидационной точности предотвратила переобучение. Дальнейшие исследования планируется сосредоточить на других трансформерах.

Для задач бинарной классификации «меланома / не меланома» и «доброкачественное/злокачественное» ViT показал точность 92,93% и 91,35% соответственно, что на 3–5% выше CNN. В трансформере были использованы частичное замораживание слоев (первые семь блоков ViT оставались заморожены, последние четыре и классификатор настраивались), скорость обучения 1е-4 с косинусным расписанием, размер батча 64, аугментация с горизонтальными отражениями и коррекцией яркости/контраста.

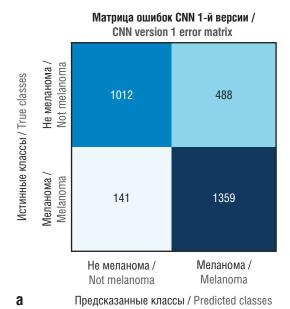
Как и в CNN, артефакты изображений (черные края, круги по центру) снижали точность модели ViT, поэтому такие данные были удалены из тренировочного набора. Пример изображений, которые были удалены, представлен на рисунке 5.

Модель на основе VIT показала значительное преимущество в сравнении с моделью CNN на основе Хсерtion в обеих задачах классификации. В задаче «меланома / не меланома» ViT достиг точности 92,9% против 88% у CNN (на 4,9% выше), при этом чувствительность составила 93,4% против 89% (+4,4%), а специфичность — 92,5% против 87% (+5,5%). В задаче «доброкачественное/злокачественное» преимущество ViT оказалось еще более выраженным: точность 91,4% против 85% (+6,4%), чувствительность 87,7% против 87% (+0,7%), а специфичность достигла впечатляющих 95% против 83% у CNN (+12%). Основные преимущества ViT заключаются в более высокой точности (в среднем на 5–6% выше) и значительно лучшей специфичности (в среднем на 5–12% выше).

Модель ViT демонстрирует превосходство в выявлении сложных паттернов благодаря способности анализировать глобальные взаимосвязи на изображении. Однако она более чувствительна к качеству входных данных. Для клинического применения, особенно в сложных случаях дифференциальной диагностики, визуальные трансформеры демонстрируют явное преимущество, но требуют более тщательной настройки.

⁸ Параметр, определяющий количество обучающих примеров, обрабатываемых за одну итерацию.





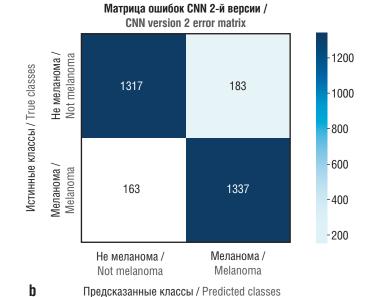




Рисунок 3. Матрицы ошибок двух версий моделей сверточных нейронных сетей (англ. convolutional neural network, CNN) MobileNetV2 и Xception (a, b) и визуального трансформера (англ. vision transformer, ViT) (c) по классификации «меланома / не меланома»

Figure 3. Error matrices of two versions of convolutional neural network (CNN) models - MobileNetV2 and Xception (a, b) and vision transformer (ViT) (c) in classifying "melanoma/non-melanoma"

Сравнение матриц ошибок двух версий моделей CNN и модели ViT по классификациям «доброкачественное/злокачественное» и «меланома / не меланома» представлены на рисунках 3, 4 и в таблице 1.

Предсказанные классы / Predicted classes

ОБСУЖДЕНИЕ / DISCUSSION

Проведенные исследования выявили критическую зависимость качества классификации от правильной подготовки тренировочной выборки. Модели, обученные с учетом оптимальной балансировки классов, продемонстрировали более высокую диагностическую точность и устойчивость классификации. В частности, для задачи «меланома / не меланома» сбалансированное распределение подклассов в классе «не меланома» позволило минимизировать ошибки классификации. Аналогично, для задачи «доброкачественное/злокачественное» равномерное распределение подклассов внутри категорий способствовало снижению количества систематических ошибок и повышению надежности моделей.

Результаты исследования показывают, что ViT превосходит CNN (MobileNetV2 и Xception) в задачах классификации кожных опухолей, достигая точности 92,93% и 91,35% против

88,47% и 85,40% для Хсерtion в задачах «меланома / не меланома» и «доброкачественное/злокачественное» соответственно. Это подтверждает способность ViT эффективно анализировать глобальные паттерны изображения, такие как асимметрия и распределение пигмента, что особенно важно для сложных случаев дерматологической диагностики.

Сравнение с другими работами подтверждает наши выводы. Например, исследование, проведенное в 2024 г., показало точность 96,15% на датасете НАМ10000 для бинарной классификации «доброкачественное/злокачественное» с использованием ViT-Google patch-32, что даже выше наших показателей – вероятно, из-за различий в предобработке данных или конфигурации модели [43]. Это подчеркивает потенциал VIT при использовании больших объемов данных, но также указывает на необходимость оптимизации предобработки для достижения максимальной производительности.

Основные ограничения ViT включают:

- необходимость больших объемов размеченных данных, что проблематично в медицине;
- высокую вычислительную сложность, увеличивающую время обучения и инференса;
- чувствительность к артефактам изображений.

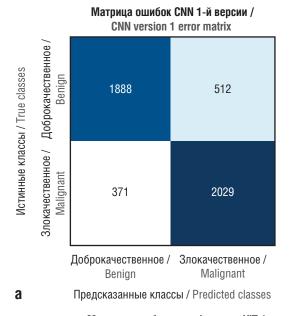






Рисунок 4. Матрицы ошибок двух версий моделей сверточных нейронных сетей (англ. convolutional neural network, CNN) MobileNetV2 и Xception (a, b) и визуального трансформера (англ. vision transformer, ViT) (c) по классификации «доброкачественное/элокачественное»

Figure 4. Error matrices of two versions of convolutional neural network (CNN) models – MobileNetV2 and Xception (**a**, **b**) and vision transformer (ViT) (**c**) in classifying "benign/malignant"

CNN также страдают от ограниченной интерпретируемости и чувствительности к сдвигам домена. Для преодоления этих проблем предлагаются:

Предсказанные классы / Predicted classes

- использование трансферного обучения и синтетических данных для ViT;
- оптимизация моделей (например, квантование для ускорения);
- улучшение предобработки для удаления артефактов;
- применение ансамблевых методов для повышения интерпретируемости и надежности.

Дальнейшие исследования должны сосредоточиться на гибридных моделях (CNN + ViT) и стандартизации протоколов валидации для медицинских данных, чтобы обеспечить их практическое внедрение.

ЗАКЛЮЧЕНИЕ / CONCLUSION

Проведенное сравнительное исследование моделей с архитектурами CNN (MobileNetV2, Xception) и модели ViT в задачах классификации кожных образований показало, что ViT демонстрирует более высокие метрики. Это связано со способностью трансформеров анализировать глобальные паттерны, такие как асимметрия и распределение пигмента. Однако ViT

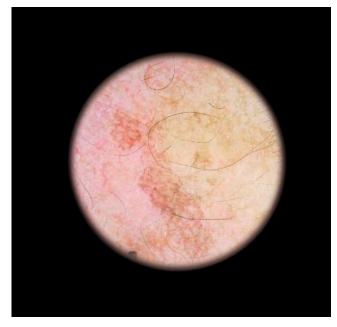


Рисунок 5. Пример изображения с артефактом **Figure 5.** Example of an image with an artifact

Таблица 1. Сравнение метрик моделей с архитектурами сверточных нейронных сетей (англ. convolutional neural network, CNN) и визуального трансформера (англ. vision transformer, ViT)

Table 1. Comparison of metrics of convolutional neural network (CNN) and vision transformer (ViT) models

Модель / Model	Архитектура / Architecture	Точность, % / Ассигасу, %	Чувствительность, % / Sensitivity, %	Специфичность, % / Specificity, %	
Меланома / не меланома // Melanoma/non-melanoma					
Визуальный трансформер / Vision transformer	ViT	92,93	93,40	92,47	
CNN 1-й версии / CNN version 1	MobileNetV2	79,03	93.40	67,47	
CNN 2-й версии / CNN version 2	Xception	88,47	89,13	87,80	
Д	оброкачественное/злока	ачественное // Benig	gn/malignant		
Визуальный трансформер / Vision transformer	ViT	91,35	87,71	95,00	
CNN 1-й версии / CNN version 1	Xception	81,60	84,54	78,67	
CNN 2-й версии / CNN version 2	Xception	85,40	87,25	83,54	

требователен к объему и качеству данных, а его эффективность снижается при наличии артефактов (например, черных краев изображений).

CNN, в первую очередь Хсерtion, остаются стабильным решением, особенно при ограниченных данных, но уступают ViT в сложных случаях дифференциальной диагностики. Ключевыми факторами, влияющими на точность обеих архитектур, стали: балансировка классов, аугментация (коррекция освещения, геометрические трансформации), настройка гиперпараметров (скорость обучения, размер батча) и очистка данных от артефактов.

Представленная работа вносит вклад в развитие методов медицинской диагностики, являясь одним из первых исследований, где проводится прямое сравнение CNN и ViT в зада-

че классификации кожных опухолей. Новизна исследования связана с тем, что архитектура модели ViT в CV стала активно применяться лишь в последние годы, и ее диагностические возможности в дерматологии изучены недостаточно.

Для клинической диагностики оптимальным решением может стать комбинация нескольких моделей (CNN + ViT) в рамках одной программы для ЭВМ (приложения). Поскольку каждая архитектура по-разному реагирует на вариации гиперпараметров, аугментации и качества входных данных, их совместное использование позволит минимизировать ошибки за счет консенсуса предсказаний. В клинической медицине, где каждый случай индивидуален, такой подход повысит надежность диагностики, особенно в сложных случаях.

ИНФОРМАЦИЯ О СТАТЬЕ	ARTICLE INFORMATION		
Поступила: 22.06.2025	Received: 22.06.2025		
В доработанном виде: 02.09.2025	Revision received: 02.09.2025		
Принята к печати: 25.09.2025	Accepted: 25.09.2025		
Опубликована: 30.09.2025	Published: 30.09.2025		
Вклад авторов	Authors' contribution		
Авторы принимали равное участие в сборе, анализе и интерпретации	The authors participated equally in the collection, analysis and interpretation		
данных. Авторы прочитали и утвердили окончательный вариант рукописи	of the data. The authors have read and approved the final version of the manuscript		
Конфликт интересов	Conflict of interests		
Авторы заявляют об отсутствии конфликта интересов	The authors declare no conflict of interests		
Финансирование	Funding		
Авторы заявляют об отсутствии финансовой поддержки	The authors declare no funding		
Этические аспекты	Ethics declarations		
Неприменимо	Not applicable		
Раскрытие данных	Data sharing		
Первичные данные могут быть предоставлены по обоснованному	Raw data could be provided upon reasonable request to the corresponding		
запросу автору, отвечающему за корреспонденцию	author		
	Publisher's note		
Комментарий издателя	Publisher's note		
Комментарий издателя Содержащиеся в этой публикации утверждения, мнения и данные были	Publisher's note The statements, opinions, and data contained in this publication were		
	1 44454		
Содержащиеся в этой публикации утверждения, мнения и данные были	The statements, opinions, and data contained in this publication were		
Содержащиеся в этой публикации утверждения, мнения и данные были созданы ее авторами, а не издательством ИРБИС (ООО «ИРБИС»).	The statements, opinions, and data contained in this publication were generated by the authors and not by IRBIS Publishing (IRBIS LLC).		
Содержащиеся в этой публикации утверждения, мнения и данные были созданы ее авторами, а не издательством ИРБИС (ООО «ИРБИС»). Издательство снимает с себя ответственность за любой ущерб, нане-	The statements, opinions, and data contained in this publication were generated by the authors and not by IRBIS Publishing (IRBIS LLC). IRBIS LLC disclaims any responsibility for any injury to people or property		
Содержащиеся в этой публикации утверждения, мнения и данные были созданы ее авторами, а не издательством ИРБИС (ООО «ИРБИС»). Издательство снимает с себя ответственность за любой ущерб, нанесенный людям или имуществу в результате использования любых	The statements, opinions, and data contained in this publication were generated by the authors and not by IRBIS Publishing (IRBIS LLC). IRBIS LLC disclaims any responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred in the		
Содержащиеся в этой публикации утверждения, мнения и данные были созданы ее авторами, а не издательством ИРБИС (ООО «ИРБИС»). Издательство снимает с себя ответственность за любой ущерб, нанесенный людям или имуществу в результате использования любых идей, методов, инструкций или препаратов, упомянутых в публикации	The statements, opinions, and data contained in this publication were generated by the authors and not by IRBIS Publishing (IRBIS LLC). IRBIS LLC disclaims any responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred in the content		
Содержащиеся в этой публикации утверждения, мнения и данные были созданы ее авторами, а не издательством ИРБИС (ООО «ИРБИС»). Издательство снимает с себя ответственность за любой ущерб, нанесенный людям или имуществу в результате использования любых идей, методов, инструкций или препаратов, упомянутых в публикации Права и полномочия	The statements, opinions, and data contained in this publication were generated by the authors and not by IRBIS Publishing (IRBIS LLC). IRBIS LLC disclaims any responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred in the content Rights and permissions		

в коммерческих целях.

ЛИТЕРАТУРА / REFERENCES

- 1. Schaumburg F., Berli C. Challenges and proposed solutions for optical reading on point-of-need testing systems. *Front Sensors*. 2023; 4. https://doi.org/10.3389/fsens.2023.1327240.
- 2. Visalini S., Kanagavalli R. A comprehensive survey of pneumonia diagnosis: image processing and deep learning advancements. In: 2023 3rd International Conference on Innovative Mechanisms for Industry Applications (ICIMIA). https://doi.org/10.1109/ICIMIA60377.2023. 10426403.
- 3. Prabha S., Gupta S., Pandey S.P. Deep learning for medical image segmentation using convolutional neural networks. In: 2024 International Conference on Optimization Computing and Wireless Communication (ICOCWC). https://doi.org/10.1109/ICOCWC60930. 2024.10470841.
- 4. Das M., Sambodhi P.P., Khare A., Naik S.A. Challenges of medical text and image processing. In: 2022 International Conference on Advancements in Smart, Secure and Intelligent Computing (ASSIC). https://doi.org/10.1109/ASSIC55218.2022.10088402.
- 5. Choudhury S., Gowri R., Babu Sena P., Dinh-Thuan D. (Eds) Intelligent Communication, Control and Devices Proceedings of ICICCD 2020: Proceedings of ICICCD 2020. https://doi.org/10.1007/978-981-16-1510-8. 6. Ламоткин А.И., Корабельников Д.И., Ламоткин И.А. и др. Искусственный интеллект в здравоохранении и медицине: история ключевых событий, его значимость для врачей, уровень развития в разных странах. ФАРМАКОЭКОНОМИКА. Современная фармакоэкономика и фармакоэпидемиология. 2024; 17 (2): 243–50. https://doi.org/10.17749/2070-4909/farmakoekonomika.2024.254.
- Lamotkin A.I., Korabelnikov D.I., Lamotkin I.A., et al. Artificial intelligence in healthcare and medicine: the history of key events, its significance for doctors, the level of development in different countries. *FARMAKOEKONOMIKA*. *Sovremennaya farmakoekonomika i farmakoepidemiologiya / FARMAKOEKONOMIKA*. *Modern Pharmacoeconomics and Pharmacoepidemiology*. 2024; 17 (2): 243–50 (in Russ.). https://doi.org/10.17749/2070-4909/farmakoekonomika.2024.254.
- 7. Ламоткин А.И., Корабельников Д.И., Ламоткин И.А. и др. Точность предварительной диагностики злокачественных меланоцитарных опухолей кожи с помощью программы искусственного интеллекта Melanoma Check. *Медицинский вестник Главного военного клинического госпиталя им. Н.Н. Бурденко.* 2025; 1: 42–51. https://doi.org/10.53652/2782-1730-2025-6-1-42-51.
- Lamotkin A.I., Korabelnikov D.I., Lamotkin I.A., et al. The accuracy of the preliminary diagnosis of malignant melanocytic skin tumors using the artificial intelligence program "Melanoma Check". *Medical Bulletin of the Main Military Clinical Hospital named after N.N. Burdenko.* 2025; 1: 42–51 (in Russ.). https://doi.org/10.53652/2782-1730-2025-6-1-42-51. 8. Zhou Z., Jin Y., Ye H., et al. Classification, detection, and segmentation performance of image-based Al in intracranial aneurysm: a systematic review. *BMC Med Imaging.* 2024; 24 (1): 164. https://doi.org/10.1186/s12880-024-01347-9.
- 9. Корабельников Д.И., Ламоткин А.И. Эффективность применения искусственного интеллекта в клинической медицине. *ФАРМАКО-ЭКОНОМИКА*. *Современная фармакоэкономика и фармакоэпидемиология*. 2025; 18 (1): 114–24. https://doi.org/10.17749/2070-4909/farmakoekonomika.2025.287.
- Korabelnikov D.I., Lamotkin A.I. The effectiveness of using artificial intelligence in clinical medicine. *FARMAKOEKONOMIKA. Sovremennaya farmakoekonomika i farmakoepidemiologiya / FARMAKOEKONOMIKA. Modern Pharmacoeconomics and Pharmacoepidemiology.* 2025; 18 (1): 114–24 (in Russ.). https://doi.org/10.17749/2070-4909/farmakoekonomika.2025.287.
- 10. Alzubaidi L., Zhang J., Humaidi A.J., et al. Review of deep learning: concepts, CNN architectures, challenges, applications, future directions. *J Big Data*. 2021; 8 (1): 53. https://doi.org/10.1186/s40537-021-00444-8. 11. LeCun Y., Bengio Y., Hinton G. Deep learning. *Nature*. 2015; 521: 436–44. https://doi.org/10.1038/nature14539.

- 12. Ламоткин А.И., Корабельников Д.И., Ламоткин И.А. Предварительная дифференциальная диагностика доброкачественных и злокачественных опухолей из эпидермальной ткани кожи с применением программы искусственного интеллекта «Derma Onko Check». Современные проблемы здравоохранения и медицинской статистики. 2025; 2: 223–42. https://doi.org/10.24412/2312-2935-2025-2-223-242.
- Lamotkin A.I., Korabelnikov D.I., Lamotkin I.A. Preliminary differential diagnosis of benign and malignant tumors from epidermal skin tissue using an artificial intelligence program "Derma Onko Check". *Current Problems of Health Care and Medical Statistics*. 2025; 2: 223–42 (in Russ.). https://doi.org/10.24412/2312-2935-2025-2-223-242.
- 13. Ламоткин А.И., Корабельников Д.И., Олисова О.Ю., Ламоткин И.А. Эффективность предварительной дифференциальной диагностики доброкачественных и злокачественных новообразований кожи с помощью программы искусственного интеллекта Derma Onko Check. *ФАРМАКОЭКОНОМИКА. Современная фармакоэкономика и фармакоэпидемиология*. 2025; 18 (2): 261–70. https://doi.org/10.17749/2070-4909/farmakoekonomika.2025.294.
- Lamotkin A.I., Korabelnikov D.I., Olisova O.Yu., Lamotkin I.A. Effectiveness of preliminary differential diagnosis of benign and malignant skin neoplasms using the Derma Onko Check artificial intelligence program. *FARMAKOEKONOMIKA. Sovremennaya farmakoekonomika i farmakoepidemiologiya / FARMAKOEKONOMIKA. Modern Pharmacoeconomics and Pharmacoepidemiology.* 2025; 18 (2): 261–70 (in Russ.). https://doi.org/10.17749/2070-4909/farmakoekonomika.2025.294.
- 14. Milletari F., Ahmadi S.A., Kroll C., et al. Hough-CNN: deep learning for segmentation of deep brain regions in MRI and ultrasound. *Computer Vision Image Underst*. 2017; 164: 92–102. https://doi.org/10.48550/arXiv.1601.07014.
- 15. Yamada M., Saito Y., Imaoka H., et al. Development of a real-time endoscopic image diagnosis support system using deep learning technology in colonoscopy. *Sci Rep.* 2019; 9 (1): 14465. https://doi.org/10.1038/s41598-019-50567-5.
- 16. Yadav D., Rathor S. Bone fracture detection and classification using deep learning approach. In: 2020 International Conference on Power Electronics & IoT Applications in Renewable Energy and its Control (PARC). https://doi.org/10.1109/PARC49193.2020.236611.
- 17. Rahman T., Chowdhury M.E., Khandakar A., et al. Transfer learning with deep convolutional neural network (CNN) for pneumonia detection using chest X-ray. *Appl Sci.* 2020; 10 (9): 3233. https://doi.org/10.3390/app10093233.
- 18. Hamamoto R., Suvarna K., Yamada M., et al. Application of artificial intelligence technology in oncology: towards the establishment of precision medicine. *Cancers*. 2020; 12 (12): 3532. https://doi.org/10.3390/cancers12123532.
- 19. Asada K., Kobayashi K., Joutard S., et al. Uncovering prognosis-related genes and pathways by multi-omics analysis in lung cancer. *Biomolecules*. 2020; 10: 524. https://doi.org/10.3390/biom10040524. 20. Kobayashi K., Bolatkan A., Shiina S., Hamamoto R. Fully-connected neural networks with reduced parameterization for predicting histological types of lung cancer from somatic mutations. *Biomolecules*. 2020; 10 (9): 1249. https://doi.org/10.3390/biom10091249.
- 21. Takahashi S., Asada K., Takasawa K., et al. Predicting deep learning based multi-omics parallel integration survival subtypes in lung cancer using reverse phase protein array data. *Biomolecules*. 2020; 10 (10): 1460. https://doi.org/10.3390/biom10101460.
- 22. Takahashi S., Sakaguchi Y., Kouno N., et al. Comparison of vision transformers and convolutional neural networks in medical image analysis: a systematic review. *J Med Syst.* 2024; 48 (1): 84. https://doi. org/10.1007/s10916-024-02105-8.
- 23. Selvaraju R.R., Cogswell M., Das A., et al. Grad-CAM: visual explanations from deep networks via gradient-based localization. In:

- 2017 Proceedings of the IEEE international conference on computer vision. https://doi.org/10.48550/arXiv.1610.02391.
- 24. Takahashi S., Takahashi M., Kinoshita M., et al. Fine-tuning approach for segmentation of gliomas in brain magnetic resonance images with a machine learning method to normalize image differences among facilities. *Cancers*. 2021; 13: 1415. https://doi.org/10.3390/cancers13061415.
- 25. Nam H., Lee H., Park J., et al. Reducing domain gap by reducing style bias. In: 2021 Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. https://doi.org/10.48550/arXiv.1910.11645.
- 26. Yan W., Wang Y., Gu S., et al. The domain shift problem of medical image segmentation and vendor-adaptation by Unet-GAN. In: Medical Image Computing and Computer Assisted Intervention—MICCAI 2019: 22nd International Conference, Shenzhen, China, October 13–17, 2019, Proceedings, Part II. https://doi.org/10.48550/arXiv.1910.13681.
- 27. Barzekar H., Patel Y., Tong L., Yu Z. MultiNet with transformers: a model for cancer diagnosis using images. arXiv:230109007. https://doi.org/10.48550/arXiv.2301.09007.
- 28. Vaswani A., Shazeer N., Parmar N., et al. Attention is all you need. In: Advances in Neural Information Processing Systems 30 (NIPS 2017). https://doi.org/10.48550/arXiv.1706.03762.
- 29. Dosovitskiy A., Beyer L., Kolesnikov A., et al. An image is worth 16×16 words: transformers for image recognition at scale. arXiv:201011929. https://doi.org/10.48550/arXiv.2010.11929.
- 30. Liu Y., Wu Y.H., Sun G., et al. Vision transformers with hierarchical attention. arXiv:210603180. https://doi.org/10.48550/arXiv.2106.03180. 31. Han K., Wang Y., Chen H., et al. A survey on vision transformer. arXiv:2012.12556. https://doi.org/10.48550/arXiv.2012.12556.
- 32. Hatamizadeh A., Yin H., Heinrich G., et al. In: 2023 Global context vision transformers. arXiv:2206.09959. https://doi.org/10.48550/arXiv.2206.09959.
- 33. He K., Gan C., Li Z., et al. Transformers in medical image analysis. *Intel Med.* 2023; 3 (1): 59–78. https://doi.org/10.1016/j.imed. 2022.07.002.

- 34. Stassin S., Corduant V., Mahmoudi S.A., Siebert X. Explainability and evaluation of vision transformers: an in-depth experimental study. *Electronics*. 2023; 13 (1): 175. https://doi.org/10.3390/electronics13010175.
- 35. Chetoui M., Akhloufi M.A. Explainable vision transformers and radiomics for COVID-19 detection in chest X-rays. *J Clin Med.* 2022; 11 (11): 3013. https://doi.org/10.3390/jcm11113013.
- 36. Dipto S.M., Reza M.T., Rahman M.N.J., et al. An XAI integrated identification system of white blood cell type using variants of vision transformer. In: Proceedings of the Second International Conference on Innovations in Computing Research (ICR'23). https://doi.org/10.1007/978-3-031-35308-6_26.
- 37. Cao Y.H., Yu H., Wu J. Training vision transformers with only 2040 images. arXiv:2201.10728. https://doi.org/10.48550/arXiv.2201.10728.
- 38. Lee S.H., Lee S., Song B.C. Vision transformer for small-size datasets. arXiv:211213492. https://doi.org/10.48550/arXiv.2112.13492. 39. Liu Y., Sangineto E., Bi W., et al. Efficient training of visual transformers with small datasets. arXiv:2106.03746. https://doi.org/10.48550/arXiv.2106.03746.
- 40. Fukushima K. Neocognitron: a self-organizing neural network model for a mechanism of pattern recognition unaffected by shift in position. *Biol Cybernetics*. 1980; 36 (4): 193–202. https://doi.org/10.1007/BF00344251.
- 41. LeCun Y., Bottou L., Bengio Y., Haffner P. Gradient-based learning applied to document recognition. *Proceedings IEEE*. 1998; 86 (11): 2278–324. https://doi.org/10.1109/5.726791.
- 42. Hamamoto R., Komatsu M., Takasawa K., et al. Epigenetics analysis and integrated analysis of multiomics data, including epigenetic data, using artificial intelligence in the era of precision medicine. *Biomolecules*. 2020; 10 (1): 62. https://doi.org/10.3390/biom10010062. 43. Himel G.M.S., Islam M.M., Al-Aff K.A., et al. Skin cancer segmentation and classification using vision transformer for automatic analysis in dermatoscopy-based noninvasive digital system. *Int J Biomed Imaging*. 2024; 2024: 3022192. https://doi.org/10.1155/2024/3022192.

Сведения об авторах / About the authors

Ламоткин Андрей Игоревич / Andrey I. Lamotkin – ORCID: https://orcid.org/0000-0001-7930-6018. eLibrary SPIN-code: 4170-7782. E-mail: lamotkin.an@mail.ru.

Корабельников Даниил Иванович, к.м.н., доцент / Daniil I. Korabelnikov, PhD, Assoc. Prof. – ORCID: https://orcid.org/0000-0002-0459-0488. eLibrary SPIN-code: 7380-7790.