

# Программные системы: теория и приложения

*Двуязычный электронный научный журнал*

№3  2024

*Bilingual Online Scientific Journal*

# Program Systems: Theory and Applications

Том 15 Выпуск 3(62) 2024 г.

## СОДЕРЖАНИЕ

*Научная статья*

МАТЕМАТИЧЕСКОЕ МОДЕЛИРОВАНИЕ

Панферов А. Д.<sup>✉</sup>, Новиков Н. А., Ульянова А. А.. *Воспроизведение отклика графена на действие внешнего электрического поля с использованием модели сильно взаимодействующих ближайших соседей* .. 3–19, 20–22

*Научная статья*

МАТЕМАТИЧЕСКОЕ МОДЕЛИРОВАНИЕ

Степанов Д. Н.<sup>✉</sup>, Тищенко И. П.. *Математическое моделирование и исследование оптимальной конфигурации оптической стереосистемы, состоящей из двух плоских зеркал* ..... 23–49, 50–53

*Научная статья*

ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ И МАШИННОЕ ОБУЧЕНИЕ

Смирнов А. В.<sup>✉</sup>. *Применение Сиамских нейронных сетей для классификации биомассы растений по визуальному состоянию* 53–72, 73–74

*Научная статья*

ИСКУССТВЕННЫЙ ИНТЕЛЛЕКТ И МАШИННОЕ ОБУЧЕНИЕ

Винокуров И. В.<sup>✉</sup>. *Восстановление текстовых последовательностей с использованием моделей глубокого обучения (Англ., Рус.)* 75–92, 93–110

*Авторский указатель* ..... 113

Click the flag at a top corner of any page to switch the language, please!

*Author index* ..... 114

*Contents* ..... 115

УДК 519.688:519.876.5

10.25209/2079-3316-2024-15-3-3-22



# Воспроизведение отклика графена на действие внешнего электрического поля с использованием модели сильно взаимодействующих ближайших соседей

Анатолий Дмитриевич **Панферов**<sup>1&2</sup>, Николай Андреевич **Новиков**<sup>2</sup>,  
Анастасия Алексеевна **Ульянова**<sup>3</sup>

<sup>1-3</sup> Саратовский государственный университет им. Н. Г. Чернышевского, Саратов, Россия

<sup>1&2</sup> [panferovad@sgu.ru](mailto:panferovad@sgu.ru)

**Аннотация.** Численное моделирование взаимодействия электромагнитного излучения с графеном позволяет воспроизводить быстро протекающие нелинейные процессы и их наблюдаемые проявления. В работе представлены результаты, полученные в процессе разработки программного решения для расчета параметров таких процессов.

Для физики графена классическим является приближение безмассовых фермионов. Однако при исследовании процессов с высокой плотностью энергии модель на основе этого приближения может оказаться за пределами своей применимости и получаемые на её основе результаты нельзя считать достоверными. Для решения этой проблемы выполнен переход к существенно более точному описанию свойств электронной подсистемы исследуемого материала, основанному на строгом учете сильного взаимодействия ближайших соседей в его кристаллической решетке.

Проведенное сравнительное тестирование двух моделей показало, что при низких энергетических характеристиках внешнего возмущения результаты совпадают. Однако, с ростом напряженности воздействующего электромагнитного поля проявляются и становятся существенными различия.

Новая точная модель имеет более сложную математическую формулировку и её использование требует больше вычислительных ресурсов. При одинаковых параметрах решаемой задачи это выражается в увеличении необходимого для выполнения расчетов времени. Относительные и абсолютные значения увеличения времени счета приведены для ряда примеров.

Полученные результаты позволяют расширить область параметров для моделирования нелинейных процессов в графене, например, генерации высокочастотных гармоник и обеспечить достоверность получаемых результатов.

**Ключевые слова и фразы:** численное моделирование, нелинейные процессы, квантовое кинетическое уравнение, модель сильно взаимодействующих ближайших соседей

**Благодарности:** Исследование выполнено за счет гранта *Российского научного фонда № 23-21-00047*<sup>ORCID</sup>

**Для цитирования:** Панферов А. Д., Новиков Н. А., Ульянова А. А. *Воспроизведение отклика графена на действие внешнего электрического поля с использованием модели сильно взаимодействующих ближайших соседей* // Программные системы: теория и приложения. 2024. Т. 15. № 3(62). С. 3–22. [https://psta.psir.ru/read/psta2024\\_3\\_3-22.pdf](https://psta.psir.ru/read/psta2024_3_3-22.pdf)

## Введение

Появление экспериментальных установок для генерации в широком диапазоне частот мощных и при этом очень коротких импульсов сделало возможным исследовать сверхбыструю динамику электронов твердых тел. В том числе двумерных кристаллов, среди которых выделяется своей необычной зонной структурой графен. Например, развивается такое новое направление, как прямое экспериментальное исследование зонной структуры различных материалов. Впервые это стало возможно при использовании фотоэлектронной спектроскопии с угловым разрешением [1]. Новыми являются применимые в условиях естественной среды чисто оптические методы: томография зонной структуры с использованием боковой полосы гармоник [2], спектроскопия высокочастотных гармоник [3, 4], и интерферометрия Блоховских электронов [5]. Для их успешного применения необходимо уметь точно моделировать нестационарную квантовую динамику электронной подсистемы материала в сильных внешних электрических полях с произвольной зависимостью от времени. Такое моделирование может реализовываться с использованием подходов, в основе которых лежит переход от рассмотрения многочастичной системы к одночастичной задаче, решаемой с использованием зависящего от времени уравнение Шредингера в той или иной форме. Точность доступных результатов будет определяться на этапе задания гамильтониана одночастичной модели. Так, применительно к графену наиболее простым и при этом эффективным вариантом является использование приближения безмассовых Дираковских фермионов (MLF – Massless Fermions) [6, 7].

Близость модели MLF к квантовой электродинамике, описывающей реальные и виртуальные электроны и позитроны, взаимодействующие посредством электромагнитного поля, позволило адаптировать для моделирования процессов в графене подход на основе квантового кинетического уравнения [8–10]. Однако приближение MLF корректно только в непосредственной окрестности точек Дирака при энергии рассматриваемых состояний не более примерно 0.5 эВ. Необходимость при моделировании мощных и высокочастотных импульсов работать с состояниями в пределах всех возможных энергий (всей зоны Бриллюэна) требует перехода к строгому учету сильного взаимодействия ближайших соседей (TBM – Tight Binding Model).

Универсальность подхода на основе квантового кинетического уравнения позволяет выполнить такой переход. В статье представлена реализация полученной модели. Выполнено сравнение предсказываемой динамики заселенности электронных состояний с результатами приближения MLF

в условиях воздействия внешнего импульсного электрического поля. В области применимости приближения MLF показана близость получаемых результатов. За её пределами демонстрируются различия в поведении моделей. На основании вычислительных экспериментов получена оценка роста ресурсоемкости задачи моделирования при использовании более точной модели.

## 1. Кинетическое уравнение для модели сильно взаимодействующих ближайших соседей

Отклик графена на действие внешнего электрического поля определяется нестационарной квантовой эволюцией его электронов, выражающейся в их перераспределении по доступным состояниям. В рамках одноэлектронного приближения совокупность таких состояний можно рассматривать как двухуровневую систему со специфической зависимостью энергии от импульса (законом дисперсии). Существуют различные подходы для описания и моделирования таких процессов. Можно непосредственно использовать уравнение Шредингера с явной зависимостью от времени [11], методы функций Грина [12, 13] или зависящего от времени функционала плотности [14, 15]. Удобным с точки зрения реализации численных методов является квантовое кинетическое уравнение [8–10]

$$(1) \quad \begin{cases} \dot{f}(\bar{p}, t) = \frac{1}{2} \lambda(\bar{p}, t) u(\bar{p}, t), \\ \dot{u}(\bar{p}, t) = \lambda(\bar{p}, t) (1 - 2f(\bar{p}, t)) - \frac{2\epsilon(\bar{p}, t)}{\hbar} v(\bar{p}, t), \\ \dot{v}(\bar{p}, t) = \frac{2\epsilon(\bar{p}, t)}{\hbar} u(\bar{p}, t), \end{cases}$$

Оно независимо определяет вероятность заселения  $f(\bar{p}, t)$  каждого из доступных состояний  $\bar{p}$  в двумерном импульсном пространстве и вспомогательные функции  $u(\bar{p}, t)$  и  $v(\bar{p}, t)$  если нам известен явный вид коэффициентов  $\lambda(\bar{p}, t)$  и  $\epsilon(\bar{p}, t)$ . Определяя состояние  $f(\bar{p}, t_{in})$ ,  $u(\bar{p}, t_{in})$  и  $v(\bar{p}, t_{in})$  в некоторый начальный момент времени  $t_{in}$  и решая задачу Коши, можно получить всю необходимую информацию о последующей эволюции. Далее начальное состояние будет определяться «вакуумными» условиями

$$(2) \quad f(\bar{p}, t_{in}) = u(\bar{p}, t_{in}) = v(\bar{p}, t_{in}) = 0.$$

В этом состоянии в отсутствие внешнего возмущающего воздействия система может находиться неограниченно долго. Поскольку нас интересует нетривиальная эволюция системы в условиях действия внешнего электрического поля,  $t_{in}$  должно предшествовать моменту его включения. Как правило, предметом рассмотрения являются короткие импульсы поля, после которых система переходит в новое стационарное состояние.

Конкретная процедура выбора используемого значения  $t_{in}$  зависит от способа задания возмущающего импульса и его характеристик.

Для получения целостной картины о поведении моделируемой системы необходимо рассматривать достаточно представительную выборку состояний (точек в импульсном пространстве) и для каждого из них решать уравнение (1). Независимость квантовой эволюции состояний обеспечивает возможность эффективного распараллеливания вычислительной процедуры. Проблематика построения выборок состояний для задач данного типа, обеспечивающих компромисс между точностью получаемых результатов и используемыми вычислительными ресурсами, рассматривалась в работах [16, 17].

Явный вид функций  $\lambda(\vec{p}, t)$  и  $\epsilon(\vec{p}, t)$  определяется гамильтонианом рассматриваемой системы и видом зависимости от времени внешнего возмущающего электрического поля. Интерес к графену с момента первых экспериментальных тестов был обусловлен свойствами его зонной структуры в окрестностях особых точек, в которых имеет место соприкосновение валентной зоны и зоны проводимости. Теория, развитая для описания наблюдаемых проявлений этих особенностей, эффективно использовала приближение безмассовых фермионов и соответствующий этой модели гамильтониан. Этот гамильтониан был использован и при адаптации метода квантового кинетического уравнения для графена. В этом случае

$$(3) \quad \lambda_{MLF}(\vec{p}, t) = eV_F^2 \frac{E_2(t)P_1 - E_1(t)P_2}{\epsilon_{MLF}^2(\vec{p}, t)},$$

$$(4) \quad \epsilon_{MLF} = V_F \sqrt{P_1^2 + P_2^2}.$$

Для обозначения этой версии коэффициентов использована аббревиатура MLF,  $e$  – элементарный заряд,  $V_F$  – скорость Ферми,  $E_1$  и  $E_2$  – компоненты внешнего электрического поля. Величины

$$(5) \quad P_1 = p_1 - eA_1(t), \quad P_2 = p_1 - eA_2(t)$$

это компоненты псевдоимпульса, соответствующего рассматриваемому состоянию,  $A_1$  и  $A_2$  – компоненты векторного потенциала внешнего электрического поля. Связь между компонентами векторного потенциала и вектора напряженности электрического поля, по определению, выражается соотношением:

$$(6) \quad E_i(t) = -\frac{\partial A_i(t)}{\partial t}, \quad i = 1, 2.$$

Однозначность определения векторного потенциала обеспечивается использованием Гамильтоновой калибровки. Приведенные выражения (3) и (4)

справедливы в системе координат с началом в одной из точек Дирака. При этом в данном приближении все такие точки эквивалентны. В ряде работ была продемонстрирована возможность использования разработанного подхода для моделирования отклика на внешние возмущения различного типа [8, 9, 18].

С ростом интереса к нелинейным явлениям в твердых телах в сильных электрических полях и появлением импульсных источников инфракрасного диапазона, способных генерировать очень короткие импульсы с высокой плотностью энергии, появилась потребность в моделировании таких процессов применительно к графену. В этом случае необходимо рассматривать и учитывать динамику состояний не только в непосредственной окрестности точек Дирака, где энергии состояний малы. Результаты, которые дает в этом случае приближение безмассовых фермионов, будут заведомо не точны.

Для решения этой проблемы необходимо использовать строгую модель ТВМ [19, 20], точно учитывающую взаимодействие с ближайшими соседями для каждого из атомов в гексагональной решетке графена. В этой модели точки Дирака формируют две неэквивалентные подрешетки и их необходимо рассматривать независимо. По этой причине удобно использовать систему отсчета с началом, расположенным между парой неэквивалентных точек Дирака. В такой системе отсчета выражения для коэффициентов системы уравнений (1) для ТВМ принимают вид [9, 21]:

$$(7) \quad \lambda_{ТВМ}(\bar{p}, t) = -\frac{4e\hbar V_F^2}{9a\varepsilon_{ТВМ}^2(\bar{p}, t)} \left( E_1(t)\sqrt{3} \left( \cos\left(\frac{\sqrt{3}aP_1}{2\hbar}\right) \cos\left(\frac{aP_2}{2\hbar}\right) + \cos\left(\frac{aP_2}{\hbar}\right) \right) + E_2(t)3 \sin\left(\frac{\sqrt{3}aP_1}{2\hbar}\right) \sin\left(\frac{aP_2}{2\hbar}\right) \right),$$

$$(8) \quad \varepsilon_{ТВМ} = \frac{2\hbar V_F}{\sqrt{3}a} \sqrt{3 - 4 \cos\left(\frac{\sqrt{3}aP_1}{2\hbar}\right) \cos\left(\frac{aP_2}{2\hbar}\right) + 2 \cos\left(\frac{aP_2}{\hbar}\right)},$$

где  $\hbar$  – приведенная постоянная Планка,  $a = 0.246$  нм – постоянная решетки графена. Здесь и далее в качестве единичного значения для  $P_1, p_1, P_2, p_2$  используется  $\hbar/a$ .

## 2. Верификация в области малых энергий

В силу того, что гамильтониан приближения MLF является результатом разложения гамильтониана ТВМ в окрестностях одной из точек Дирака (в которых энергия состояний принимает нулевое значение) с сохранением

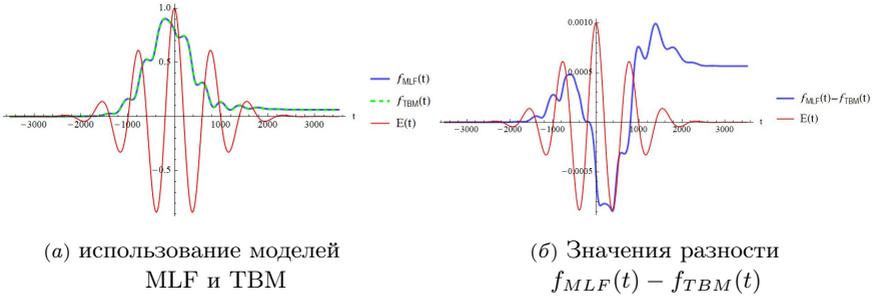
только линейных по  $p_1$  и  $p_2$  членов [20], результаты для версий (3), (4) и (7), (8) при моделировании ограниченных по энергии возбуждаемых состояний возмущений должны совпадать.

Энергия, которую внешняя электромагнитная волна может передать электронам материала, зависит от её частоты и амплитуды. Так, при частоте  $\nu = 5 \times 10^{12}$  Гц и амплитуде электрического поля  $E_0 = 3 \times 10^5$  В/м энергия генерируемых возбуждений не будет превышать примерно  $2 \times 10^{-2}$  эВ. Рассмотрим короткий импульс с такими параметрами, определяя его через компоненты векторного потенциала в форме

$$(9) \quad \begin{aligned} A_1(t) &= -\frac{E_0}{2\pi\nu} \sin(2\pi\nu t) \exp\left(-\frac{t^2}{2\tau^2}\right), \\ A_2(t) &= 0. \end{aligned}$$

Это линейно поляризованный импульс с гауссовой огибающей с длительностью  $\tau$ , которую определим условием  $2\pi\nu\tau = 6$ .

Результаты воспроизведения поведения функции распределения для состояния с  $p_1 = 0.0$ ,  $p_2 = 2.0975$  в окрестностях точки Дирака с энергией 0.0083 эВ в условиях действия импульса внешнего поля (9) представлены на рисунке 1а. На всех графиках в статье масштабы для напряженности электрического поля  $E(t)$  и единица шкалы времени условны.



(а) использование моделей  
MLF и TBM

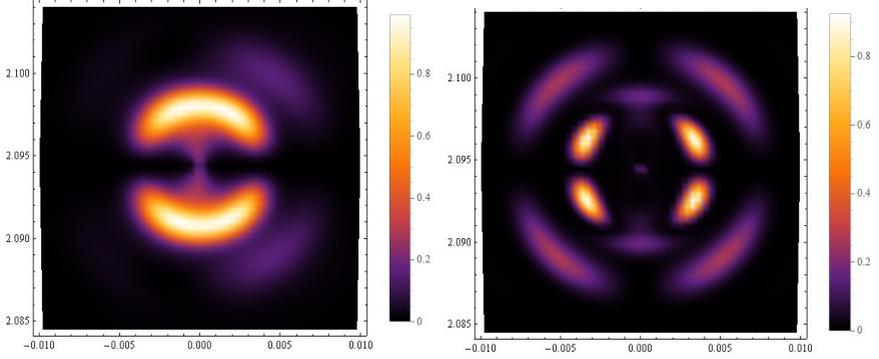
(б) Значения разности  
 $f_{MLF}(t) - f_{TBM}(t)$

Рисунок 1. Результаты воспроизведения эволюции функций распределения для состояния в окрестности точки Дирака с энергией 0.0083 эВ;

Для наглядности функции распределения показаны вместе с зависящим от времени полем возмущающего импульса. В масштабе рисунка  $f_{MLF}(t)$  и  $f_{TBM}(t)$  совпадают на всем интервале значений времени. Однако между ними нет строгого равенства, значения и зависимость от времени разности  $f_{MLF}(t) - f_{TBM}(t)$  показано на рисунке 1б. Она относительно невелика, но наблюдаема и воспроизводима при использовавшейся точности

моделирования. Так, в конечном стационарном состоянии после завершения действия поля  $f_{MLF}(t_{end}) = 0.0618699$  и  $f_{TBM}(t_{end}) = 0.0613054$ .

Обобщенная картина  $f(\vec{p})$  в двумерном импульсном пространстве (пространстве состояний) приведена на рисунке 2а для момента времени  $t = 0$ , когда напряженность внешнего поля максимальна, а на рисунке 2б – для момента времени  $t_{end}$ , когда действие возмущения завершено.



(а) в момент максимума внешнего поля при  $t = 0$

(б) после завершения действия импульса поля

Рисунок 2. Функции распределения для области состояний  $-0.01 < p_1 < 0.01, 2.084 < p_2 < 2.104$

В силу близости и визуальной неразличимости результатов для моделей MLF и TBM приведены результаты для первой модели. Различия показаны на рисунке 3.

Они малы и локализованы в окрестностях точки Дирака (центр показанной области). На этом рисунке проявляется характерное отличие TBM модели – её анизотропия. Хотя наложение симметрии распределения возбуждений в линейно поляризованном поле смазывает картину, на рисунках различимо проявление трехлучевой симметрии этой модели.

### 3. Выход за границы применимости приближения безмассовых фермионов

Для сравнения моделей в условиях возбуждения состояний с более высокими уровнями энергии увеличим напряженность электрического поля в десять раз до  $E_0 = 3 \times 10^6$  В/м не меняя другие параметры импульса. При таких параметрах приближение безмассовых фермионов должно еще

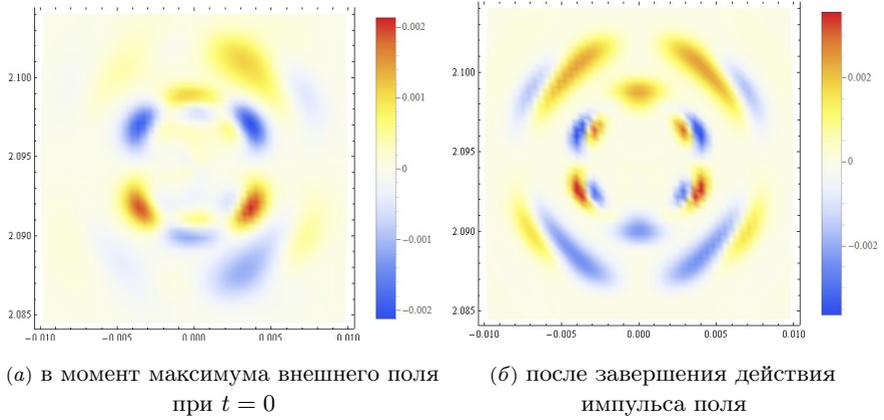


Рисунок 3. Распределение различий в результатах использования моделей  $f_{MLF}(\vec{p}, t) - f_{TBM}(\vec{p}, t)$  для области состояний  $-0.01 < p_1 < 0.01, 2.084 < p_2 < 2.104$

достаточно точно описывать реакцию материала. Действительно, результаты моделирования функции распределения для моментов максимума внешнего поля при  $t = 0$  и после завершения действия импульса визуально не различимы.

Они представлены на рисунке 4. Показанная здесь область импульсного

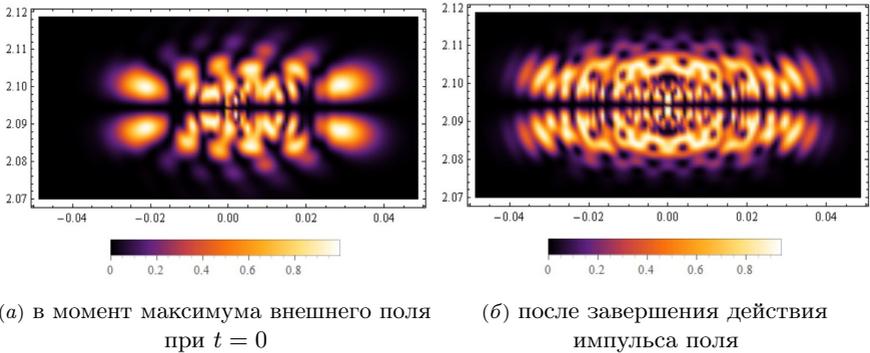
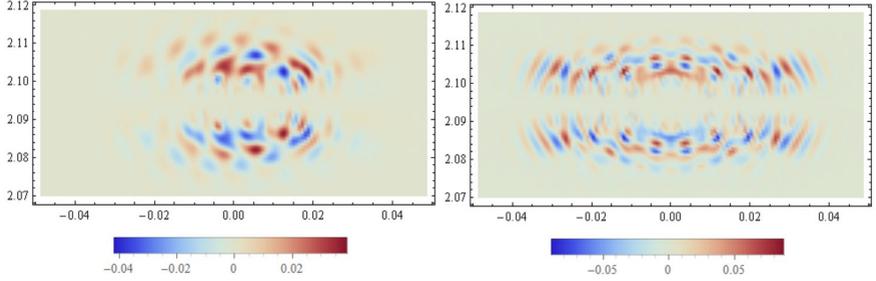


Рисунок 4. Функция распределения при  $E_0 = 3 \times 10^6$  В/м для области состояний  $-0.05 < p_1 < 0.05, 2.069 < p_2 < 2.119$ .

пространства больше чем на рисунке 3 по вертикальной оси в 2.5 раза, а по горизонтальной оси в 5 раз. Наблюдается распространение возбуждений далее от точки Дирака, расположенной в центре. Растянutosть вдоль

горизонтальной оси соответствует направлению действия внешнего поля.

Различия результатов для рассматриваемых моделей показаны на рисунке 5.

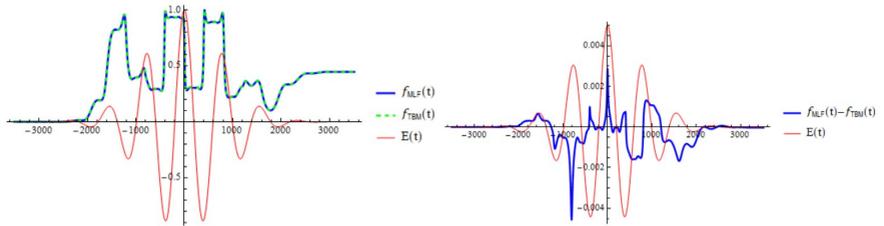


(а) в момент максимума внешнего поля при  $t = 0$       (б) после завершения действия импульса поля

Рисунок 5. Распределение различий между  $f_{MLF}(\vec{p}, t)$  и  $f_{TBM}(\vec{p}, t)$  в окрестности точки Дирака при  $E_0 = 3 \times 10^6$  В/м для области состояний  $-0.05 < p_1 < 0.05, 2.069 < p_2 < 2.119$ .

Распределение разностных значений ассоциировано с распределением максимальных значений функций  $f_{MLF}(\vec{p}, t)$  и  $f_{TBM}(\vec{p}, t)$ . При этом диапазон наблюдаемых разностных значений  $f_{MLF}(\vec{p}, t) - f_{TBM}(\vec{p}, t)$  достигает почти 0.1 в то время как в предыдущем случае он не превышал 0.006. Меняется и характер поведения функций распределения во времени.

Для демонстрации этого на рисунке 6 представлены результаты для той же точки в импульсном пространстве с энергией состояния 0.0083 эВ, что и на рисунке 1.



(а) функций распределения      (б) разности  $f_{MLF}(t) - f_{TBM}(t)$

Рисунок 6. Поведение в окрестности точки Дирака с энергией 0.0083 эВ при  $E_0 = 3 \times 10^6$  В/м

Зависимость от времени в окрестностях максимума импульса приобретает ступенчатый характер с резкими скачками значений, совпадающими по времени с максимумами напряженности поля. Увеличение амплитуды поля еще в десять раз до  $E_0 = 3 \times 10^7$  В/м выводит за границы применимости модели MLF. Получающийся результат для функции распределения приведен на рисунке 7.

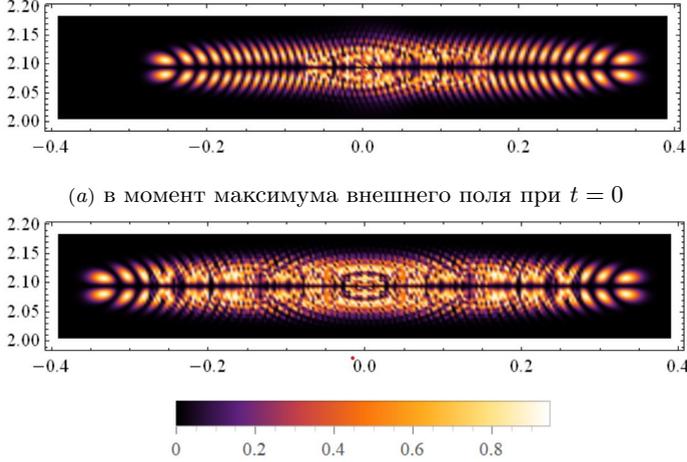
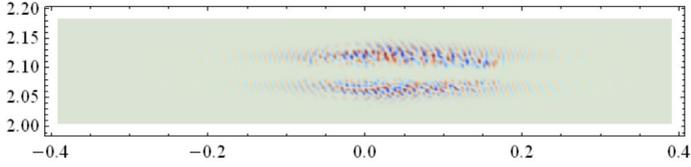


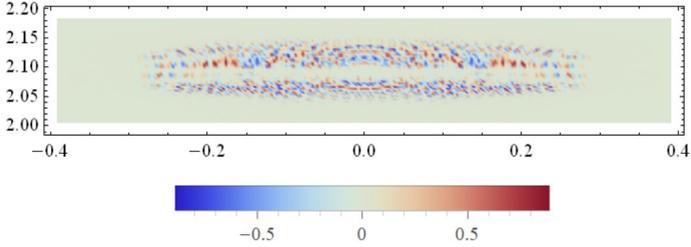
РИСУНОК 7. Функция распределения  $f_{MLF}(\bar{p}, t)$  при  $E_0 = 3 \times 10^7$  В/м для области состояний  $-0.09 < p_1 < 0.09$ ,  $1.869 < p_2 < 2.319$ .

Теперь в процесс вовлечены состояния далекие от точки Дирака. Как и в предыдущих случаях приведены распределения только для  $f_{MLF}(\bar{p}, t)$ . Хотя различия между  $f_{MLF}(\bar{p}, t)$  и  $f_{TBM}(\bar{p}, t)$  существенны, в силу сложной структуры распределения возбужденных состояний и доминирования симметрии обусловленной направлением поля, их визуальное выявление затруднительно и приведено в форме распределения разностных значений на рисунке 8.

Теперь диапазон  $-1 \lesssim f_{MLF}(\bar{p}, t) - f_{TBM}(\bar{p}, t) \lesssim 1$  охватывает все допустимые значения с учетом того, что, по определению, для фермионов  $0 \leq f(\bar{p}, t) \leq 1$ . Особенности поведения  $f_{MLF}(\bar{p}, t)$  и  $f_{TBM}(\bar{p}, t)$  во времени в условиях действия сильного внешнего поля проиллюстрированы



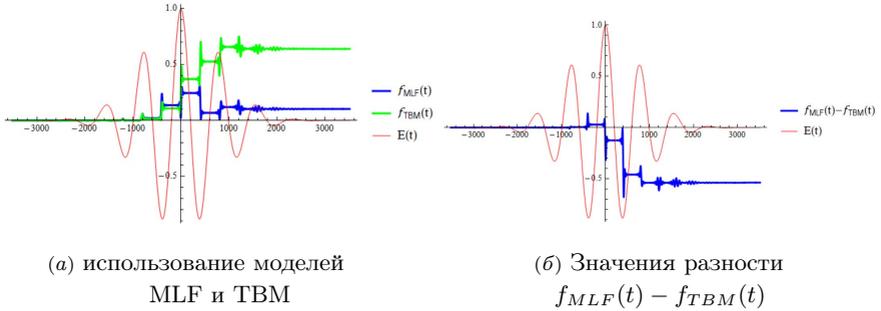
(а) в момент максимума внешнего поля при  $t = 0$



(б) после завершения действия импульса поля

Рисунок 8. Распределение различий между  $f_{MLF}(\bar{p}, t)$  и  $f_{TBM}(\bar{p}, t)$  при  $E_0 = 3 \times 10^7$  В/м для области состояний  $-0.09 < p_1 < 0.09, 1.869 < p_2 < 2.319$ .

на примере состояния с большой разностью конечных значений в точке  $p_1 = 0.0, p_2 = 2.138$  на рисунке 9.



(а) использование моделей MLF и TBM

(б) Значения разности  $f_{MLF}(t) - f_{TBM}(t)$

Рисунок 9. Поведение функций распределения для состояния с  $p_1 = 0.0, p_2 = 2.138$  при  $E_0 = 3 \times 10^7$  В/м

Отмеченные для рисунка 6 особенности теперь проявляются более явно. Каждый текущий максимум внешнего поля вызывает осцилляции

и последующий резкий переход к новому, относительно стабильному значению. Это свидетельствует о жестком поведении системы уравнений при используемом наборе параметров.

#### 4. Варианты реализации и требования к ресурсам

При переходе к новой модели (7), (8) необходимо оценить возможные изменения в потребных вычислительных ресурсах. Поскольку изменения заключаются только в переопределении аналитического представления коэффициентов системы уравнений (1) при прочих равных условиях время решения аналогичных задач не должно существенно измениться.

В таблице 1 приведены значения времени выполнения заданий на тестовой системе с характеристиками: процессор 11th Gen *Intel*<sup>®</sup>*Core*<sup>®</sup> i5-11400 @ 2.60GHz (шесть физических ядер), оперативная память 32.0 Гб.

Таблица 1. Результаты тестов на время выполнения моделирования

	тест 1	тест 2	тест 3	tes 4
Число состояний	24616	51184	67900	24616
Время для модели MLF, с.	96.895	297.85	8042.8	401.39
Время для модели TBM, с.	106.70	337.85	8970.9	447.91
Увеличение времени счета	10.12%	13.43%	11.54%	11.59%

Система функционирует под управлением ОС Windows 11 Pro. Моделирование выполнялось в версии системы, реализованной средствами Wolfram Mathematica. Обеспечивалось независимое решение системы уравнений (1) для различных состояний на шести физических ядрах процедурой *NDSolve* с установленными параметрами точности решения *AccuracyGoal*  $\rightarrow$  10 и *PrecisionGoal*  $\rightarrow$  10.

Физические параметры моделируемого воздействия соответствуют представлявшимся выше. Так, тест 1 выполнялся для амплитуды электрического поля  $E_0 = 3 \times 10^5$  В/м, тест 2 для  $E_0 = 3 \times 10^6$  В/м, тест 3 и тест 4 для  $E_0 = 3 \times 10^7$  В/м. Частота во всех случаях оставалась постоянной  $\nu = 5 \times 10^{12}$  Гц. Для каждого набора параметров формировалась собственная сетка состояний по процедуре, являющейся модификацией

процедуры, описанной в [17]. Поэтому число моделируемых состояний в тестах различно. Увеличение их числа коррелировано с усложнением наблюдаемой картины распределения возбужденных состояний. Последний тест выполнялся при параметрах внешнего поля, тех же что и в третьем случае, но на сетке состояний из первого теста.

Из представленных данных следует, что в рассматриваемом случае усложнение модели приводит к увеличению времени счета не более чем на 14%. Однако необходимо отметить, что на время счета влияют и физические параметры моделируемого процесса. При сравнении результатов тестов 1–3 видно, что время выполнения задания растет существенно быстрее роста количества состояний в используемой сетке. Выше, при сравнении рисунка 1, рисунка 6 и рисунка 9 на качественном уровне отмечалось усложнение поведения  $f(\bar{p}, t)$  во времени. Именно это является причиной роста времени на воспроизведение поведения функции распределения в каждом состоянии при наложенных фиксированных требованиях к точности решения.

Для явной демонстрации этого факта был выполнен тест с номером 4. В нем использован тот же набор состояний, что и в тесте 1, но напряженность внешнего поля больше на два порядка. Это приводит к более чем четырехкратному росту времени счета. Причем и для старой, и для новой модели.

Высокие требования решаемой задачи к вычислительным ресурсам делают необходимым использование высокопроизводительных кластерных систем. Поэтому вычислительная процедура была реализована на языке C++ средствами библиотеки Boost с распределением воспроизведения эволюции рассматриваемых состояний по параллельным процессам средствами MPI. Влияние выбора модели на время счета для этой реализации было протестировано на примере линейно поляризованного импульса со следующим набором параметров: частота  $\nu = 80 \times 10^{12}$  Гц, амплитуда электрического поля  $E_0 = 1.6 \times 10^9$  В/м и  $2\pi\nu t = 12$ . Использовались сетки с числом состояний от 1180 до 75672. При запуске задачи с распараллеливанием на 40 процессов на вычислительном узле с двумя CPU Intel® Xeon® Gold 6338 @ 2.00GHz (32 физических ядра, гипертрейдинг), оперативная память 768.0 Гб были получены следующие результаты, представленные в таблице 2.

При запуске теста для счета 75672 состояний с распараллеливанием на 128 процессов с использованием двух нод приведенной конфигурации

Таблица 2. Результаты для кластерной реализации

Число состояний	1180	4728	18916	75672
Время для модели MLF, с.	10.108	40.702	124.53	501.86
Время для модели ТВМ, с.	11.784	50.003	156.66	616.54
Увеличение времени счета	16.58%	22.85%	25.80%	22.85%

увеличение времени счета для новой модели составило 24.18%, что близко к результатам, показанным в таблице 2.

### Заключение

В работе представлены результаты тестирования модели для описания поведения графена в условиях действия на него высокочастотных электромагнитных импульсов с высокой плотностью энергии. Основой является кинетическое уравнение, обеспечивающее возможность численного воспроизведения нестационарной квантовой эволюции электронов рассматриваемой физической системы. Явный вид уравнения определен с точным учетом параметров взаимодействия ближайших соседей в кристаллической решетке материала. В отличие от использовавшегося ранее приближения безмассовых фермионов, ограниченного в применимости процессами с малой энергией возбуждаемых электронных состояний, эта модель свободна от такого ограничения и может быть использована для исследования нелинейных процессов в более широком диапазоне параметров.

Сравнительное тестирование показало совпадение результатов с приближением безмассовых фермионов в условиях, удовлетворяющих требованиям применимости обеих моделей. Если же параметры моделируемого воздействия не удовлетворяют таким требованиям, предсказываемые результаты квантовой эволюции электронных состояний становятся различны вплоть до полного несовпадения. Последнее подтверждает корректность мотивации о необходимости перехода к новой модели для получения в этих условиях правильных результатов.

Представленная модель первоначально была реализована средствами Wolfram Mathematica и протестирована в многопоточном режиме на рабочей станции. Кластерная версия реализации выполнена на C++

с использованием МРІ и библиотеки Boost. Результаты её тестирования также представлены.

Разработанные решения предоставляют новые возможности для моделирования поведения графена в условиях действия на него электромагнитного излучения высокой интенсивности.

### Список использованных источников

- [1] Zhang H., Pincelli T., Jozwiak Ch., Kondo T., Ernstorfer R., Sato T., Zhou S. *Angle-resolved photoemission spectroscopy* // Nature Reviews Methods Primers.– 2022.– Vol. **2**.– id. 54.– 22 pp. [doi](#) ↑4
- [2] Mikhailov S. A. *Non-linear electromagnetic response of graphene* // Europhysics Letters.– 2007.– Vol. **79**.– id. 27002.– 5 pp. [doi](#) ↑4
- [3] Ishikawa K. L. *Nonlinear optical response of graphene in time domain* // Phys. Rev. B.– 2010.– Vol. **82**.– id. 201402. [doi](#) ↑4
- [4] Yoshikawa N. *High-harmonic generation in graphene enhanced by elliptically polarized light excitation* // Science.– 2017.– Vol. **356**.– No. 6339.– Pp. 736–738. [doi](#) ↑4
- [5] Cha S., Kim M., Kim Y., Choi Sh., Kang S., Kim H., Yoon S., Moon G., Kim T., Lee Y. W., Cho G. Y., Park M. J., Kim Ch-J., Kim B. J., Lee JD., Jo M-H., Kim J. *Gate-tunable quantum pathways of high harmonic generation in graphene* // Nature Communication.– 2022.– Vol. **13**.– id. 6630.– 10 pp. [doi](#) ↑4
- [6] Novoselo K. S., Geim A. K., Morozov S. V., Jiang D., Katsnelson M. I., Grigorieva I. V., Dubonos S. V., Firsov A. A. *Two-dimensional gas of massless Dirac fermions in graphene* // Nature.– 2005.– Vol. **438**.– Pp. 197–200. [doi](#) ↑4
- [7] Castro Neto A. H., Guinea F., Peres N. M. R., Novoselov K. S., Geim A. K. *The electronic properties of graphene* // Rev. Mod. Phys.– 2009.– Vol. **81**.– No. 1.– id. 109. [doi](#) ↑4
- [8] Panferov A., Smolyansky S., Blaschke D., Gevorgyan N. *Comparing two different descriptions of the I-V characteristic of graphene: theory and experiment*, XXIV International Baldin Seminar on High Energy Physics Problems “Relativistic Nuclear Physics and Quantum Chromodynamics” (Baldin ISHEPP XXIV) // EPJ Web Conf.– 2019.– Vol. **204**.– id. 06008.– 6 pp. [doi](#) ↑4, 5, 7
- [9] Smolyansky S., Panferov A., Blaschke D., Gevorgyan N. *Nonperturbative kinetic description of electron-hole excitations in graphene in a time dependent electric field of arbitrary polarization* // Particles.– 2019.– Vol. **2**.– No. 2.– Pp. 208–230. [doi](#) ↑4, 5, 7
- [10] Smolyansky S. A., Blaschke D. B., Dmitriev V. V., Panferov A. D., Gevorgyan N. T. *Kinetic equation approach to graphene in strong external fields* // Particles.– 2020.– Vol. **3**.– No. 2.– Pp. 456–476. [doi](#) ↑4, 5

- [11] Boolakee T., Heide Ch., Wagner F., Ott Ch., Schlecht M., Ristein J., Weber H., Hommelhoff P. *Length-dependence of light-induced currents in graphene* // J. Phys. B: At. Mol. Opt. Phys.– 2020.– Vol. **53**.– No. 15.– id. 154001.– 5 pp. doi ↑5
- [12] Ke M., Asmar M. M., Tse W. K. *Nonequilibrium RKKY interaction in irradiated graphene* // Physical Review Research.– 2020.– Vol. **2**.– No. 3.– id. 033228. doi ↑5
- [13] Li J., Han J. E. *Nonequilibrium excitations and transport of Dirac electrons in electric-field-driven graphene* // Phys. Rev. B.– 2018.– Vol. **97**.– No. 20.– id. 205412. doi ↑5
- [14] Chen Zi-Yu., Qin R. *Circularly polarized extreme ultraviolet high harmonic generation in graphene* // Optics Express.– 2019.– Vol. **27**.– No. 3.– Pp. 3761–3770. doi ↑5
- [15] Li P., Shi R., Lin P., Ren X. *First-principles calculations of plasmon excitations in graphene, silicene, and germanene* // Phys. Rev. B.– 2023.– Vol. **107**.– No. 3.– id. 035433. doi ↑5
- [16] Панферов А. Д., Новиков Н. А., Трунов А. А. *Моделирование поведения графена во внешних электрических полях* // Программные системы: теория и приложения.– 2021.– Т. **12**.– № 1(38).– С. 3–19. URL doi ↑6
- [17] Панферов А. Д., Поснова Н. В., Ульянова А. А. *Моделирование поведения двухуровневой квантовой системы с использованием масштабируемых регулярных сеток* // Программные системы: теория и приложения.– 2023.– Т. **14**.– № 2(57).– С. 27–47. URL doi ↑6, 15
- [18] Панферов А. Д., Новиков Н. А. *Характеристики индуцированного излучения в условиях действия на графен коротких высокочастотных импульсов* // Известия Саратовского университета. Новая серия. Серия: Физика.– 2023.– Т. **23**.– № 3.– С. 254–264. doi ↑7
- [19] Reich S., Maultzsch J., Thomsen C., Ordejon P. *Tight-binding description of graphene* // Phys. Rev. B.– 2002.– Vol. **66**.– No. 3.– id. 035412. doi ↑7
- [20] Katsnelson M. I. *The Physics of Graphene*, 2nd ed.– Cambridge University Press.– 2020.– ISBN 9781108617567.– id. 425. doi ↑7, 8
- [21] Панферов А. Д., Щербаков И. А. *Реализация квантового кинетического уравнения для графена на основе модели сильного взаимодействия ближайших соседей* // Известия Саратовского университета. Новая серия. Серия: Физика.– 2024.– Т. **24**.– № 3.– С. 198–208. doi ↑7

Поступила в редакцию	01.04.2024;
одобрена после рецензирования	29.05.2024;
принята к публикации	09.08.2024;
опубликована онлайн	10.09.2024.

Рекомендовал к публикации

к.т.н. Д. В. Зубов

## Информация об авторах:



### Анатолий Дмитриевич Панферов

к.ф.-м.н., зам. начальника УЦИТ Саратовского государственного университета им. Н.Г. Чернышевского. Научные интересы: высокопроизводительные вычисления, параллельное программирование, численное решение квантовых кинетических уравнений, моделирование процессов вакуумного рождения частиц в КЭД, генерации носителей в полупроводниках в том числе бесщелевых, процессов на ранних стадиях столкновения релятивистских ядер



0000-0003-2332-0982

*e-mail:*



### Николай Андреевич Новиков

Саратовский государственный университет им. Н.Г. Чернышевского. Научные интересы: моделирование физических процессов на высокопроизводительных вычислительных системах, параллельное программирование.



0000-0003-1259-1867

*e-mail:*



### Анастасия Алексеевна Ульянова

Саратовский государственный университет им. Н.Г. Чернышевского. Научные интересы: моделирование физических процессов на высокопроизводительных вычислительных системах, параллельное программирование.



0000-0001-9519-9822

*e-mail:*

*Авторы внесли равный вклад в подготовку публикации.*

*Декларация об отсутствии личной заинтересованности: благополучие авторов не зависит от результатов исследования.*



# Simulation the response of graphene to an external electric field using the exact tight-binding model

Anatolii Panferov<sup>1</sup>, Nikolay Novikov<sup>2</sup>, Anastasiya Ulyanova<sup>3</sup>

1-3 Saratov State University, Saratov, Russia

 panferovad@sgu.ru

**Abstract.** Numerical simulation of the interaction of electromagnetic radiation with graphene allows us to reproduce fast nonlinear processes and their observed manifestations. The paper presents the results obtained in the process of developing a software solution for calculating the observed parameters of such processes.

In graphene physics, the massless fermion approximation is classical. However, in the study of processes with high energy density, model based on this approximation are beyond the limits of their applicability and the results obtained on their basis can not be considered reliable. To solve this problem, a transition to a substantially more accurate model based on a strict account of the nearest-neighbor interaction in the crystal lattice (tight-binding model) has been made.

Comparative testing of these two models shows that at low energy characteristics of the external perturbation the results coincide. However, as the energy characteristics of the affecting electromagnetic field increase, the divergence of the results becomes apparent and grows.

The new exact model has a more complex mathematical formulation and requires more computational resources. When using the same hardware configuration it is expressed in the increase of counting time. Relative and absolute values for a number of examples are given.

The obtained results allow us to expand the range of parameters for modeling of nonlinear processes in the considered material, for example, generation of high-frequency harmonics and ensure its reliability. (*In Russian*).

**Key words and phrases:** numerical simulation, nonlinear processes, quantum kinetic equation, tight-binding model

2020 *Mathematics Subject Classification:* 65Z05; 81-04; 81T40

**Acknowledgments:** The study was supported by the Russian Science Foundation grant No. 23-21-00047, <https://rscf.ru/project/23-21-00047/>

**For citation:** Anatolii Panferov, Nikolay Novikov, Anastasiya Ulyanova. *Simulation the response of graphene to an external electric field using the exact tight-binding model*. Program Systems: Theory and Applications, 2024, **15**:3(62), pp. 3–22. (*In Russ.*). [https://psta.psiras.ru/read/psta2024\\_3\\_3-22.pdf](https://psta.psiras.ru/read/psta2024_3_3-22.pdf)

## References

- [1] H. Zhang, T. Pincelli, Ch. Jozwiak, T. Kondo, R. Ernstorfer, T. Sato, S. Zhou. “Angle-resolved photoemission spectroscopy”, *Nature Reviews Methods Primers*, **2** (2022), id. 54, 22 pp. [doi](#)
- [2] S. A. Mikhailov. “Non-linear electromagnetic response of graphene”, *Europhysics Letters*, **79** (2007), id. 27002, 5 pp. [doi](#)
- [3] K. L. Ishikawa. “Nonlinear optical response of graphene in time domain”, *Phys. Rev. B*, **82** (2010), id. 201402. [doi](#)
- [4] N. Yoshikawa. “High-harmonic generation in graphene enhanced by elliptically polarized light excitation”, *Science*, **356**:6339 (2017), pp. 736–738. [doi](#)
- [5] S. Cha, M. Kim, Y. Kim, Sh. Choi, S. Kang, H. Kim, S. Yoon, G. Moon, T. Kim, Y. W. Lee, G. Y. Cho, M. J. Park, Ch-J. Kim, B. J. Kim, JD. Lee, M-H. Jo, J. Kim. “Gate-tunable quantum pathways of high harmonic generation in graphene”, *Nature Communication*, **13** (2022), id. 6630, 10 pp. [doi](#)
- [6] K. S. Novoselo, A. K. Geim, S. V. Morozov, D. Jiang, M. I. Katsnelson, I. V. Grigorieva, S. V. Dubonos, A. A. Firsov. “Two-dimensional gas of massless Dirac fermions in graphene”, *Nature*, **438** (2005), pp. 197–200. [doi](#)
- [7] Castro Neto A. H., F. Guinea, N. M. R. Peres, K. S. Novoselov, A. K. Geim. “The electronic properties of graphene”, *Rev. Mod. Phys.*, **81**:1 (2009), id. 109. [doi](#)
- [8] A. Panferov, S. Smolyansky, D. Blaschke, N. Gevorgyan. “Comparing two different descriptions of the I-V characteristic of graphene: theory and experiment”, XXIV International Baldin Seminar on High Energy Physics Problems “Relativistic Nuclear Physics and Quantum Chromodynamics” (Baldin ISHEPP XXIV), *EPJ Web Conf.*, **204** (2019), id. 06008, 6 pp. [doi](#)
- [9] S. Smolyansky, A. Panferov, D. Blaschke, N. Gevorgyan. “Nonperturbative kinetic description of electron-hole excitations in graphene in a time dependent electric field of arbitrary polarization”, *Particles*, **2**:2 (2019), pp. 208–230. [doi](#)
- [10] S. A. Smolyansky, D. B. Blaschke, V. V. Dmitriev, A. D. Panferov, N. T. Gevorgyan. “Kinetic equation approach to graphene in strong external fields”, *Particles*, **3**:2 (2020), pp. 456–476. [doi](#)
- [11] T. Boolakee, Ch. Heide, F. Wagner, Ch. Ott, M. Schlecht, J. Ristein, H. Weber, P. Hommelhoff. “Length-dependence of light-induced currents in graphene”, *J. Phys. B: At. Mol. Opt. Phys.*, **53**:15 (2020), id. 154001, 5 pp. [doi](#)
- [12] M. Ke, M. M. Asmar, W. K. Tse. “Nonequilibrium RKKY interaction in irradiated graphene”, *Physical Review Research*, **2**:3 (2020), id. 033228. [doi](#)
- [13] J. Li, J. E. Han. “Nonequilibrium excitations and transport of Dirac electrons in electric-field-driven graphene”, *Phys. Rev. B*, **97**:20 (2018), id. 205412. [doi](#)
- [14] Zi-Yu. Chen, R. Qin. “Circularly polarized extreme ultraviolet high harmonic generation in graphene”, *Optics Express*, **27**:3 (2019), pp. 3761–3770. [doi](#)
- [15] P. Li, R. Shi, P. Lin, X. Ren. “First-principles calculations of plasmon excitations in graphene, silicene, and germanene”, *Phys. Rev. B*, **107**:3 (2023), id. 035433. [doi](#)

- [16] A. D. Panferov, N. A. Novikov, A. A. Trunov. “Simulate the behavior of graphene in external electric fields”, *Program Systems: Theory and Applications*, **12**:1(48) (2021), pp. 3–19 (in Russian).  
- [17] A. D. Panferov, N. V. Posnova, A. A. Ulyanova. “Simulation the behavior of a two-level quantum system using scalable regular grids”, *Program Systems: Theory and Applications*, **14**:2(57) (2023), pp. 27–47 (in Russian).  
- [18] A. D. Panferov, N. A. Novikov. “Characteristics of induced radiation under the action of short high-frequency pulses on graphene”, *Izv. Sarat. Univ. Physics*, **23**:3 (2023), pp. 254–264 (in Russian).  
- [19] S. Reich, J. Maultzsch, C. Thomsen, P. Ordejon. “Tight-binding description of graphene”, *Phys. Rev. B*, **66**:3 (2002), id. 035412. 
- [20] M. I. Katsnelson. *The Physics of Graphene*, 2nd ed., Cambridge University Press, 2020, ISBN 9781108617567, id. 425. 
- [21] A. D. Panferov, I. A. Scherbakov. “Tight-binding implementation of the quantum kinetic equation for graphene”, *Izvestiya of Saratov University. Physics*, 2024, no. 3, pp. 198–208. 

УДК 004.942:535.318

 10.25209/2079-3316-2024-15-3-23-53

## Математическое моделирование и исследование оптимальной конфигурации оптической стереосистемы, состоящей из двух плоских зеркал

Дмитрий Николаевич Степанов<sup>1✉</sup>, Игорь Петрович Тищенко<sup>2</sup>

<sup>1,2</sup>Институт программных систем им. А. К. Айламазяна РАН, Вельсково, Россия

<sup>1✉</sup>[mitek1989@mail.ru](mailto:mitek1989@mail.ru)

**Аннотация.** Статья посвящена математическому моделированию и оптимизации конфигурации оптической стереосистемы, состоящей из видеокamеры и двух плоских зеркал. Отличие данного исследования от ранее проведенных — учет большого количества ограничений на конфигурацию оптической системы: величина стереобазы, размеры зеркал, общие габариты оптической системы, отсутствие двойного отражения световых лучей, недопущение ситуации, когда видеокamera отражается в зеркалах. Выполнена постановка задачи условной оптимизации для поиска оптимальной конфигурации рассматриваемой оптической системы. В качестве целевой функции выбран периметр прямоугольника, ограничивающего габариты оптической системы. Численное решение задачи было найдено с использованием пакета SciPy. Полученные результаты расширяют теорию компьютерного зрения и могут быть использованы в создании и исследовании систем компьютерного зрения для робототехнических комплексов.

**Ключевые слова и фразы:** машинное зрение, оптические приборы, математическое моделирование, стереозрение, оптимизация, катоптрическая система

Для цитирования: Степанов Д. Н., Тищенко И. П. *Математическое моделирование и исследование оптимальной конфигурации оптической стереосистемы, состоящей из двух плоских зеркал* // Программные системы: теория и приложения. 2024. Т. 15. № 3(62). С. 23–53.  
[https://psta.psir.ru/read/psta2024\\_3\\_23-53.pdf](https://psta.psir.ru/read/psta2024_3_23-53.pdf)

## Введение

Системы технического зрения широко применяются для решения задач, связанных с неразрушающим контролем качества деталей и материалов, создания трехмерных моделей реальных объектов, в биометрических системах распознавания, в автономных роботах и др. Основными инструментами для восстановления 3D-структуры наблюдаемой сцены в настоящее время являются:

- лидары — источники лазерного излучения; принцип действия основан на измерении времени движения отраженных лазерных лучей. Схожий принцип действия имеют времяпролетные (ToF-) камеры.
- источники структурированной подсветки: специальный излучатель (проектор) проецирует на сцену некоторый паттерн известной конфигурации, а видеокамера выполняет съемку сцены, подсвеченной проектором. Детектирование элементов паттерна на снимке с видеокамеры позволяет восстановить 3D-координаты точек, на которые попадает паттерн.
- съемка сцены несколькими фото- или видеокамерами с разных ракурсов. Для неподвижных объектов также используется вариант, когда снимки с разных ракурсов выполняются одной камерой.

Первые лидары и времяпролетные камеры были достаточно дорогими, громоздкими и ресурсоемкими устройствами, но сейчас в свободной продаже есть много доступных и компактных устройств подобных классов, а их технические характеристики позволяют использовать их для решения многих практических задач. Их важнейшее преимущество перед решениями на основе видеокамер состоит в том, что лидары и ToF-камеры позволяют получить данные о 3D-структуре окружающей обстановки без дополнительных нетривиальных алгоритмов калибровки сенсоров, а также алгоритмов обработки и анализа сенсорных данных. Работа систем на основе нескольких видеокамер основана на решении задачи стереосопоставления (*stereo correspondence problem*), в которой необходимо находить соответствия между пикселями изображений, выполненных с разных ракурсов и/или в разные моменты времени. Данная задача является одной из самых сложных в компьютерном зрении и не имеет универсального алгоритма ее решения. Если же используется структурированная подсветка, то необходимость в решении задачи стереосопоставления отпадает, но появляется задача поиска элементов подсветки на снимках с видеокамеры. А она тоже может быть нетривиальной, особенно в условиях неконтролируемого освещения.

Для восстановления трехмерной структуры сцены также используются катадиоптрические системы, в состав которых входят видеокамера

и комбинация преломляющих и отражающих элементов. На одной светочувствительной матрице создается сразу несколько изображений. Это приводит к тому, что каждая виртуальная видеокамера будет иметь меньшие углы обзора, чем исходная видеокамера, а проблема стереосопоставления остается актуальной, но для получения стереоснимков достаточно одной видеокамеры. Следовательно, в задаче калибровки оптической системы уменьшается количество неизвестных, поскольку внутренние параметры виртуальных камер идентичны. Далее будут описаны примеры практического применения подобных систем, в которых они имеют преимущества перед решениями на основе лидаров и ToF-камер.

Из катадиоптрических систем можно выделить диоптрические (используются только преломляющие элементы) и катоптрические (только отражающие элементы). В качестве примеров работ по исследованию систем из первой группы можно указать статью [1]: исследуются вопросы калибровки системы, в которой эффект стерео достигается за счет использования призмы перед камерой. В статье [2] описана система, в которой похожая призма располагается внутри камеры, между системой линз и светочувствительной матрицей, что позволяет на одном изображении получать два снимка с разных ракурсов.

Катоптрические системы можно условно разделить по типу используемых зеркал, плоских или криволинейных. По сравнению с использованием плоских зеркал, криволинейные зеркала позволяют добиться более широких углов обзора, но видимые размеры объектов становятся меньше, а геометрические искажения на изображениях — более заметными, особенно по краям зеркал. Схожие особенности имеют системы, использующие камеры с широкоугольными объективами типа «рыбий глаз» (англ. fish-eye).

Катоптрические системы с плоскими зеркалами можно классифицировать по количеству используемых зеркал. Однозеркальные системы (пример — статья [3]) наиболее простые, но имеют узкую область применения ввиду слишком малого угла обзора. Один из способов решения данной проблемы — вращение зеркала: например, в работе [4] зеркало вращается в одной плоскости, а ось вращения зеркала совпадает с направлением оптической оси камеры. В работе [5] представлена система панорамного видения, основанная на использовании камеры и плоского зеркала, вращающегося вокруг двух осей. В статье [6] описана похожая система всенаправленного панорамного видения, в которую входит зеркало, вращающееся со скоростью несколько десятков оборотов в секунду, а также специальный алгоритм управления затвором камеры.

В качестве примеров работ по трехзеркальным системам можно привести статьи [3, 7–9]: например, в статье [3] исследуется оптимальное

расположение трех зеркал для получения ректифицированной стереопары с помощью одной камеры. В статьях [7, 8] описана система из одной видеокамеры, трех обычных зеркал и светоделиителя (англ. beamsplitter), который половину света пропускает, а вторую половину — отражает. Изображения двух виртуальных камер получаются наложенными друг на друга (т.е. нуждаются в постобработке). Приведены уравнения для расчета размера зеркал, уравнения для восстановления 3D-координат точки, наблюдаемой на снимке со стереоустановки, а также алгоритм построения карты диспаратитетов на основе преобразования Фурье. В статье [9] для трехзеркальной системы выведены уравнения эпиполярных ограничений (уравнения, которые связывают координаты пикселей, соответствующих одному и тому же объекту на изображениях одной сцены с разных ракурсов).

Принцип работы трехзеркальных систем, как правило, сводится к тому, что на одну часть светочувствительной матрицы проецируются лучи, претерпевшие отражение от одного зеркала, а на другую часть матрицы — отражение от двух других зеркал. На таком принципе основана и оптическая система, описанная в статье [10]: система из 9 зеркал организована в виде двух пирамид, расположенных друг напротив друга. Поле обзора камеры делится на несколько непересекающихся регионов, и в каждом из них световой поток делится на две вышеуказанные части. Выполнен расчет оптимальной конструкции такой системы, представлен способ ее калибровки. Выполнен расчет погрешностей в оценивании 3D-координат наблюдаемой точки при известных погрешностях при измерении проекций этой точки на паре снимков с двух виртуальных камер. Аналогичная оптическая система исследована и в статье [11].

В работе [12] описана система из четырех плоских зеркал, угол между двумя внутренними зеркалами равен  $90^\circ$ . Выведено уравнение для расчета стереобазы между двумя виртуальными камерами. Рассчитана конфигурация системы для обеспечения нужного угла обзора. В статье [13] рассмотрена почти такая же система, но два внешних зеркала могут менять свою ориентацию. Выполнен расчет стереобазы, расчет размеров и расположения зеркал, расчет конфигурации области пересечения полей зрения двух виртуальных камер. Как и в работе [10], произведен расчет погрешностей в оценивании 3D-координат наблюдаемой точки. Статья [14] также посвящена расчету оптимальной конфигурации четырехзеркальной системы, с учетом желаемого угла обзора, рабочей дистанции, точности измерений и общих физических размеров оптической системы. В работе [15] выполнено сравнение возможностей четырехзеркальной системы и стереосистемы из двух видеокамер. Примером статей, в которых описано

практическое применение четырехзеркальных систем, является работа [16], посвященная построению 3D-моделей лопастей авиационных двигателей, а также статья [17], в которой описано оригинальная стереонасадка из четырех зеркал, которая позволяет получать стереоснимки с помощью обычного смартфона.

В работе [18] описан оригинальный подход для получения панорамных стереоизображения с широким углом обзора с помощью одной камеры и стереонасадки, состоящей из множества плоских зеркал. Панорамным стереоснимкам посвящена и статья [19]: несколько видеокamer наблюдают конструкцию из набора зеркал, составляющих грани пирамиды.

Данная работа посвящена двухзеркальным системам. На одном фотоприемнике формируется сразу два изображения, чаще всего одно из них занимает левую половину снимка, а второе — правую. Подобным системам посвящено достаточно много работ, но нужно заметить, что в некоторых исследованиях рассматриваются частные случаи конфигурации оптической системы. Например, в статье [20] выведены уравнения эпиллярных ограничений для двухзеркальной системы, но принято допущение, что одно из зеркал располагается параллельно плоскости изображения камеры. Одним из авторов предыдущей статьи, Shree K. Nayar, в работе [21] предложена математическая модель формирования изображений в системах из одной камеры и одного зеркала (плоского, конического, сферического, параболического, эллиптического или гиперболического). Выполнена аналитическая оценка пространственного разрешения подобной оптической системы, а также оценка уровня размытия, вызванного расфокусировкой. В системе, описанной в статье [22], зеркала располагаются симметрично относительно оптической оси камеры и соприкасаются друг с другом, камера при этом наблюдает отражения объектов, располагающиеся за камерой. Выполнен расчет стереобазы с использованием угла между зеркалами и расстояния до зеркал, а также оценка оптимального взаимного расположения зеркал, чтобы в них не было видно друг друга, и чтобы обе виртуальные камеры могли наблюдать точку, располагающуюся на бесконечности. Показано, что при уменьшении угла между зеркалами увеличивается угол обзора и одновременно увеличивается стереобаза.

Но более общие случаи конфигураций также рассматриваются: например, авторы статей [23, 24] построили математическую модель двухзеркальной системы, вывели уравнения эпиллярных ограничений, предложили алгоритм вычисления фокальной длины камеры по изображениям со стереосистемы, а также рассчитали угол обзора стереосистемы. В работе [25] описана система, состоящая из двух плоских зеркал и RGB-D камеры, которая совмещает в себе обычную видеокamerу и средство для

построения карт глубин. Два зеркала позволяют следить за обстановкой в двух направлениях: перед мобильным роботом и за ним. Оптическая ось камеры при этом направлена в точку соединения зеркал. В работе [26] также предложена математическая модель двухзеркальной системы, предложен алгоритм вычисления угла между зеркалами, а также алгоритм вычисления положения и ориентации камеры относительно зеркал по двум наблюдаемым точкам (положение вычисляется с точностью до некоторого положительного коэффициента). В статье [27] двухзеркальная оптическая система используется для измерения уровня вибрации на поверхности устройства для воспроизведения звука (колонка), причем правая часть изображения формируется из световых лучей, которые претерпели отражение от обоих зеркал, а левая часть — лучами, которые не отражались от зеркал. Работа [28] посвящена расчету конфигурации калейдоскопической стереосистемы, состоящей из одной камеры и двух или трех плоских зеркал. В статье [29] предложен способ калибровки подобной калейдоскопической системы. В работе [30] два плоских зеркала используются совместно с дихроичным фильтром, одна из сторон которых пропускает световые волны, соответствующие красному цвету, а другая сторона отражает волны, соответствующие синему.

Что касается практического применения катоптрических систем (в частности, двухзеркальных), то можно выделить следующее:

- мобильные наземные роботы, решающие задачи SLAM и/или детектирования и распознавания целевых объектов окружающей обстановки [25]. В работе [31] исследовалась стереонасадка, которая разрабатывалась для применения на мобильных роботах. Конечно, здесь следует заметить, что использование системы зеркал в данной задаче возможно только в тех случаях, когда требуемый размер стереобазы относительно невелик, иначе зеркальная система окажется слишком громоздкой. В частности, это несколько ограничивает использование подобных систем на беспилотных летательных аппаратах.
- промышленная робототехника: например, в статье [26] описывается робот-манипулятор, в рабочем пространстве которого располагаются два или более плоских зеркала. Робот оснащен видеокамерой, зеркала позволяют повысить точность работы манипулятора за счет наблюдения рабочего пространства с разных ракурсов. Отдельно взятый лидар или ToF-камера могут обеспечить построение 3D-модели объекта только с одного ракурса.
- изменение уровня вибрации на рабочих поверхностях приборов для воспроизведения звука [27]. Одним из способов решения данной задачи является стереозрение, которое имеет ряд преимуществ перед

другими методами: возможность проводить измерения по всему рабочему полю и низкая чувствительность к окружающему акустическому шуму. Высокоскоростные камеры, которые используются для решения при решении данной задачи, являются дорогим специализированным оборудованием, и поэтому использование зеркал для получения эффекта стерео может радикально уменьшить стоимость всей оптической системы. Лидары и ToF-камеры имеют значительно меньшее быстродействие (FPS, Frames per Second — количество кадров в секунду), чем высокоскоростные видеокамеры.

- получение трехмерных моделей различных объектов (получение их цифровых копий): в работе [29] использование видеокамеры и двух плоских зеркал позволило получать пять изображений одного объекта с разных ракурсов.

Наконец, можно выделить еще одно преимущество решений на основе видеокамер по сравнению с лидарами и ToF-камерами. Видеокамеры имеют более высокое разрешение и позволяют получить информацию о внешнем виде объекта (например, его текстуру, цвет, надписи на нем), лидары и ToF-камеры позволяют получить только пространственную информацию. Если стоит задача детектирования и распознавания целевых объектов на сенсорных данных в видимом диапазоне, то лидара или ToF-камеры будет недостаточно.

Обзор показал, что для расчета оптимальной конфигурации подобных зеркальных систем используется небольшое количество параметров: в частности, для двухзеркальных — только угол обзора и размер стереобазы. Кроме того, в существующих моделях игнорируется наличие практических ограничений на размер и конфигурацию отдельных элементов подобных оптических систем (это замечание относится и к моделям с большим количеством зеркал). Например, слишком большие зеркала может быть или невозможно изготовить, или невозможно установить из-за специфики того технического средства, на котором функционирует катоптрическая система (например, летающий или едущий автономный робот могут иметь ограничения на максимальный размер полезной нагрузки). Кроме того, в существующих моделях не затрагивается проблема многократного отражения: если в одном из зеркал отражается другое зеркало, то результирующее изображение будет содержать неинформативные области. Наконец, сама видеокамера может наблюдать зеркала — этой ситуации также следует избегать, но в разработанных ранее моделях подобная проблема игнорируется. Таким образом, актуальной видится разработка такой математической модели двухзеркальной катоптрической системы, в которой бы учитывались перечисленные практические аспекты проектирования подобных систем.

В конечном итоге, предлагается выполнить формальную постановку задачи условной оптимизации, которая бы учитывала все эти ограничения, для поиска оптимальной конструкции оптической системы. В данной работе будут использованы некоторые выкладки из статьи [32], которая посвящена разработке математической модели двухзеркальной системы, предложенная модель отличается учетом дисторсии на изображениях с реальных видеокамер.

Существуют различные программы для расчета и оптимизации оптических систем, одной из наиболее популярных является Zemax. Рассматриваемый класс задач имеет важную особенность: требуется, чтобы световой поток от каждого из пары зеркал в итоге проецировался только на одну из половин площади фотоприемника (плоскости изображения), а два зеркала в совокупности должны полностью захватывать поле зрения камеры, т.е. камера не должна наблюдать ничего, помимо зеркал. Это достигается не только определенным положением и ориентацией каждого из зеркал, но и их определенными размерами. Zemax подобный режим не поддерживает: в оптической схеме перед фотоприемником может быть только один оптический элемент, световые лучи от которого поступают на фотоприемник. Кроме того, в Zemax напрямую нельзя задать размеры оптического элемента «Плоское зеркало» (только его положение и угол наклона). Более того, в изученных статьях из списка литературы не было найдено упоминаний о том, что для исследования оптической системы был использован Zemax или какой-то аналог, причем это справедливо и для самых свежих публикаций.

## 1. Математическая модель оптической стереосистемы из двух плоских зеркал

На рисунке 1 приведена иллюстрация модели оптической стереосистемы из двух плоских зеркал, подробное ее описание доступно в работе [32]. Используется вид сверху, начало координат (точка  $O$ ) расположено в оптическом центре камеры, ось  $OX$  направлена вправо от камеры, ось  $OY$  — вниз (на рисунке 1 она направлена от наблюдателя, используется вид сверху), ось  $OZ$  — вперед (оптическая ось). Центр объектива камеры совпадает с началом координат, на рисунке 1 объектив изображен в виде тонкой линзы. Фокальная длина камеры равна  $f$ , плоскости обоих зеркал перпендикулярны плоскости  $OXZ$ . Отрезки  $AB_1$  и  $B_2C$  соответствуют плоскостям зеркал, вместе они перекрывают все поле зрения камеры. Точки  $B_1$  и  $B_2$  располагаются на оси  $OZ$ . Два световых луча, проходящих через точку  $P$  пространства, отражаются от зеркал и проецируются на матричный фотоприемник  $F_1F_2$  (плоскость изображения камеры)

в точки  $q_1$  и  $q_2$  соответственно (углы падения равны углам отражения). Точки  $P_1$  и  $P_2$  — наблюдаемые образы точки  $P$ . Отрезки  $AB_1$  и  $B_2C$  лежат на прямых с уравнениями  $z = \tan(\phi_1)x + b_1$  и  $z = \tan(\phi_2)x + b_2$ , здесь  $\phi_1$  и  $\phi_2$  — углы между прямыми и положительным направлением оси  $OX$ ,  $\phi_{1,2} \in [0, \pi]$ .

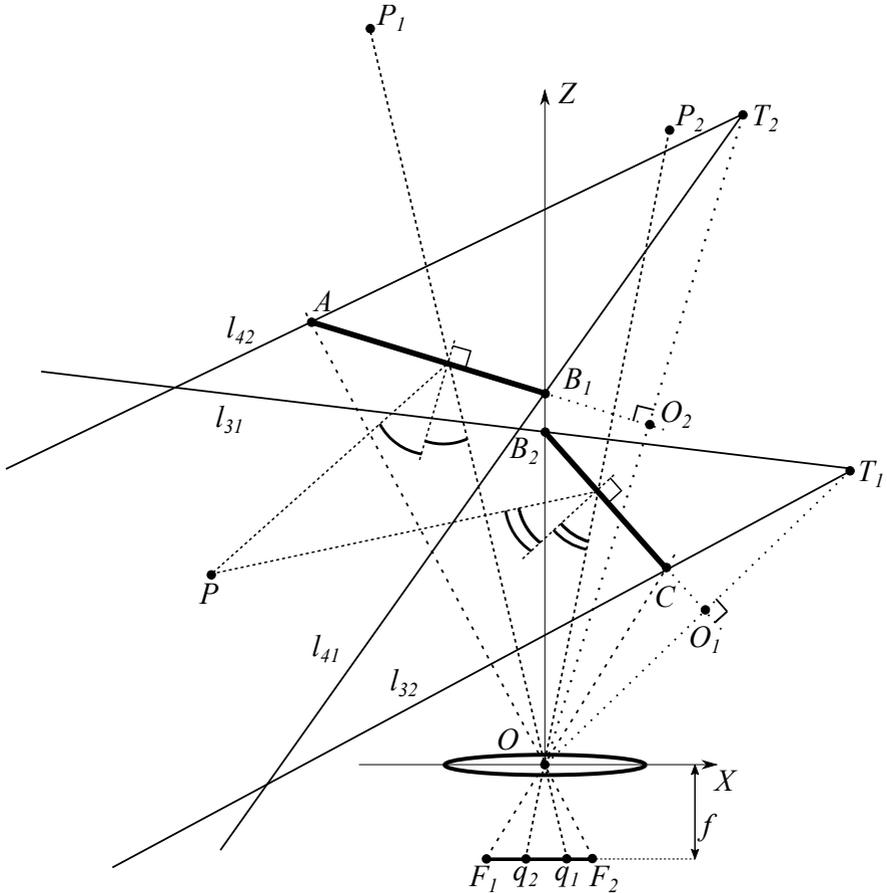


Рисунок 1. Модель оптической стереосистемы из видеокамеры и двух плоских зеркал

Введем и обоснуем использование двух допущений в используемой модели. Первое допущение: параметры используемой видеокамеры (размер матрицы, фокусное расстояние, разрешение матрицы, количество линз, их взаимное положение и конфигурация) являются неизменными и не входят

в состав параметров, которые требуется оптимизировать. Единственный параметр, который характеризует непосредственно камеру в рассматриваемой задаче — это ее фиксированный угол обзора (Field of View — FOV). В рассмотренных статьях не было обнаружено моделей, в которых бы оптические характеристики фото- или видеокамер являлись оптимизируемыми и изменяемыми параметрами, т.е. камера жестко задана. Авторам статей достаточно только угла обзора, что позволяет смоделировать камеру с помощью модели «камера-обскура» и перспективной проекции. Также угол обзора может быть вычислен при известном разрешении камеры и фокальной длине, последняя может быть вычислена в ходе калибровки либо с использованием известного фокусного расстояния и размера одного элемента светочувствительной матрицы. Чувствительность, отношение сигнал/шум, размер матрицы, параметры линз, апертура и др. — все эти параметры камер в статьях не рассматриваются. Создание собственной фото/видеокамеры — намного более трудная задача, по сравнению с созданием системы из зеркальных элементов и/или разделителей лучей (beam splitter) и готовой цифровой камеры. Таким образом, допущение о фиксированных характеристиках фото/видеокамер вполне отвечает действительности, с которой сталкиваются другие исследователи.

Второе допущение — дисторсия на снимках либо незначительна, либо устранена программными методами (камера при этом должна быть заранее откалибрована). Если же учитывать в математической модели коэффициенты дисторсии, то вывод аналитических зависимостей будет невозможен из-за невозможности получить аналитическое решение системы нелинейных уравнений, которую порождает математическая модель. Современные цифровые камеры, как правило, обладают незначительной дисторсией. Исключение — камеры с большим или сверхбольшим углом обзора, но для таких камер используются другие математические модели (отличные от модели «камера-обскура» и перспективной проекции), в настоящем исследовании подобные камеры не рассматриваются. В рассмотренных статьях других исследователей не было обнаружено, чтобы дисторсия линз принималась во внимание.

Сделаем еще несколько допущений:

- зеркала примыкают друг к другу, т.е.,  $B_1 = B_2 = B$  и  $b_1 = b_2 = b$ . Такое допущение присутствует в части рассмотренных статей, также оно имеет место в той оптической системе, которая далее будет исследована с помощью предложенной математической модели.
- ось  $OZ$  делит угол обзора пополам:  $\angle AOB = \angle COB = \alpha$ . Это соответствует ситуации, когда положение оптического центра камеры незначительно отличается от центра светочувствительной матрицы. Для большинства современных камер это условие выполняется.

- выполняется следующий набор условий:

$$(1) \quad \alpha \in \left(0, \frac{\pi}{2}\right), \phi_1 > \frac{\pi}{2} + \alpha, \phi_1 > \phi_2, \phi_{1,2} \in \left(\frac{\pi}{2}, \pi\right].$$

Условие  $\alpha \geq \frac{\pi}{2}$  может выполняться только для сверхширокоугольных камер, которые в данном исследовании не рассматриваются. Если  $\phi_1 \leq \frac{\pi}{2} + \alpha$ , то прямая  $AB$  не будет пересекаться с лучом  $OA$ , то есть, точка  $A$  не будет лежать на левой границе поля зрения камеры. Если  $\phi_1 \leq \phi_2$ , то поля зрения виртуальных камер не будут пересекаться (данная закономерность была выведена эмпирическим путем, с помощью геометрического построения различных конфигураций рассматриваемой оптической системы). Если  $\phi_{1,2} \leq \frac{\pi}{2}$ , то камера будет наблюдать заднюю поверхность зеркала.

Поскольку использование зеркал делает изображения зеркально повернутыми относительно вертикальной оси, то необходимо развернуть изображения еще раз. Пусть  $\check{P}_1$  и  $\check{P}_2$  — трехмерные координаты наблюдаемых образов точки  $P$  на изображении с оптической стереосистемы, которое было подвергнуто зеркальному повороту. В статье [32] были выведены следующие формулы для вычисления координат  $\check{P}_{1,2}$ :

$$(2) \quad \check{P}_{1,2} = \begin{bmatrix} \check{X}_{1,2} \\ \check{Y}_{1,2} \\ \check{Z}_{1,2} \end{bmatrix} = \check{R}_{1,2}P + \check{T}_{1,2} = R_{1,2}^t (P - T_{1,2}),$$

$$R_{1,2} = \begin{bmatrix} -\cos 2\phi_{2,1} & 0 & \sin 2\phi_{2,1} \\ 0 & 1 & 0 \\ -\sin 2\phi_{2,1} & 0 & -\cos 2\phi_{2,1} \end{bmatrix}, T_{1,2} = b \begin{bmatrix} -\sin 2\phi_{2,1} \\ 0 \\ 1 + \cos 2\phi_{2,1} \end{bmatrix},$$

$$\check{P}_2 = \widehat{R}^t (\check{P}_1 - \widehat{T}), \widehat{R} = R_1^t R_2, \widehat{T} = R_1^t (T_2 - T_1).$$

Матрицы  $R_1$  и  $R_2$  задают ориентации систем координат, привязанных к правой и левой виртуальным камерам соответственно (обозначим эти системы как СК1 и СК2), относительно системы координат, привязанной к реальной монокулярной камере. Верхний индекс  $t$  — операция транспонирования.  $T_1$  и  $T_2$  — координаты оптических центров виртуальных камер, начала СК1 и СК2. Пара  $(\widehat{R}, \widehat{T})$  задает ориентацию и положение СК2 относительно СК1.

Пусть  $(x_i, y_i)$  — двумерные координаты проекции точки  $P$  на изображении с  $i$ -той виртуальной камеры,  $M_i$  — матрица внутренних параметров  $i$ -той камеры (camera intrinsic matrix), строение матриц  $M_i$  в контексте данной задачи изложено в статье [32]. 3D-координаты точки  $P$  неизвестны,

как и координаты  $\check{P}_1$  и  $\check{P}_2$ . Матрицы  $M_i$  связывают координаты  $\check{P}_1$  и  $\check{P}_2$  с 2D-координатами проекций точки  $P$ :

$$(3) \quad w_i \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = M_i \check{P}_i \Rightarrow \check{P}_i = w_i M_i^{-1} \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix}.$$

Здесь  $w_1$  и  $w_2$  — неизвестные масштабирующие множители. Выражение  $\check{P}_2 = \widehat{R}^t (\check{P}_1 - \widehat{T})$  можно переписать следующим образом:

$$(4) \quad w_2 M_2^{-1} \begin{bmatrix} x_2 \\ y_2 \\ 1 \end{bmatrix} = \widehat{R}^t \left( w_1 M_1^{-1} \begin{bmatrix} x_1 \\ y_1 \\ 1 \end{bmatrix} - \widehat{T} \right).$$

Нетрудно убедиться, что это выражение является переопределенной системой из трех линейных алгебраических уравнений относительно двух неизвестных  $w_1$  и  $w_2$ , которые могут быть найдены, к примеру, с помощью метода наименьших квадратов. Необходимо только заметить, что дисторсия на снимке с видеокамеры должна быть предварительно устранена.

Используя координаты точек  $T_{1,2}$ , нетрудно доказать, что они являются отражением точки  $O$  относительно двух зеркал (см. рисунок 1,  $O_1 \in B_2C, OT_1 \perp B_2C, O_1O = O_1T_1$ ,  $O_2 \in B_1A, OT_2 \perp B_1A, O_2O = O_2T_2$ ). Угол обзора правой виртуальной камеры задается лучами  $l_{41}$  (лежит на отрезке  $T_2B_1$ ) и  $l_{42}$  (лежит на отрезке  $T_2A$ ). Аналогично задается угол обзора левой виртуальной камеры: лучи  $l_{31}$  и  $l_{32}$ . Также нетрудно доказать, что углы обзора обеих виртуальных камер равны  $\alpha$ : поскольку выполняются равенства  $\triangle OB_2O_1 = \triangle T_1B_2O_1$  и  $\triangle OCO_1 = \triangle T_1CO_1$  (по первому признаку равенства треугольников), то  $\angle B_2T_1C = \angle B_2OC = \alpha$ . Доказательство по аналогичной схеме используется для случая, когда точка  $O_1$  располагается внутри отрезка  $B_2C$ , а также для левого зеркала (отрезок  $AB_1$ ).

## 2. Ограничения, связанные со величиной стереобазы

Расстояние между камерами (в данном случае виртуальными) — один из основных параметров, характеризующих оптическую стереосистему:

$$(5) \quad \begin{aligned} T_1T_2 &= 2b |\sin(\phi_1 - \phi_2)|, \\ \phi_1 > \phi_2 &\Rightarrow T_1T_2 = 2b \sin(\phi_1 - \phi_2). \end{aligned}$$

Если система из двух камер откалибрована (известны их внутренние параметры и пара  $(\hat{R}, \hat{T})$ ), то к полученной стереопаре можно применить процедуру геометрического выравнивания — ректификации [33]. Изображения в итоге получаются такими, как если бы они были выполнены двумя камерами с одинаковыми внутренними параметрами и с одинаковой ориентацией ( $\hat{R}$  становится равной единичной матрице), а смещение между их оптическими центрами имеет место только вдоль оси  $OX$ , т.е.  $\hat{T} = [T \ 0 \ 0]^t$ . При таких условиях, коэффициенты увеличения обеих камер станут одинаковыми, а проекция любой точки пространства, которая наблюдаема на обоих снимках, будет располагаться в пикселях, имеющих одинаковый номер строки, что значительно облегчает решение задачи стереосопоставления. Пусть проекция этой точки располагается на левом и правом снимках в пикселях с номерами столбцов  $x_l$  и  $x_r$  соответственно, величина  $d = x_l - x_r$  именуется диспаратетом. Тогда расстояние до точки рассчитывается по формуле  $Z = Z(d) = \frac{fT}{d}$  [33]. Оценим погрешность  $\Delta Z$  в вычислении функции  $Z(d)$ , если известна некоторая априорная оценка  $\Delta d$  погрешности в вычислении диспаратета (она зависит от особенностей используемого алгоритма поиска стереосоответствия):

$$(6) \quad \Delta Z = \left| \frac{\partial Z}{\partial d} \right| \Delta d = \left( \frac{f^2 T^2}{d^2} \right) \frac{1}{fT} \Delta d = \frac{Z^2}{fT} \Delta d.$$

Рассмотрим следующую задачу: известно  $Z_0$  — расстояние до наблюдаемого объекта, необходимо найти параметры оптической системы, при которых погрешность в вычислении расстояния составляет не более  $\Delta Z_0$ :

$$(7) \quad \frac{Z_0^2}{fT} \Delta d \leq \Delta Z_0 \Rightarrow T \geq T_{min} = \frac{Z_0^2 \cdot \Delta d}{f \cdot \Delta Z_0}.$$

Это выражение имеет следующий смысл: какой должна быть минимальная величина стереобазы, чтобы при известной фокальной длине  $f$  обеих камер, известной погрешности  $\Delta d$  в вычислении диспаратета, погрешность в вычислении расстояния до объектов с помощью триангуляции не превосходила величину  $\Delta Z_0$  (если про объект известно, что он находится на расстоянии не более чем  $Z_0$ ). Например, если автономный подвижный робот решает задачу навигации внутри помещений, то величина  $Z_0$  может быть равна, скажем, 3-4 метра. Априорная оценка этой величины зависит от скорости движения робота, быстродействия алгоритмов навигации и технических характеристик камеры (к примеру, достаточно удаленные объекты будут, скорее всего, плохо различимы на снимках, но в то

же время, в силу удаленности в данный момент они не представляют опасности для робота).

В итоге ограничение на минимальную величину стереобазы задается следующим образом:

$$(8) \quad \text{base}(\phi_1, \phi_2, b) := 2b \sin(\phi_1 - \phi_2) \geq T_{min}.$$

### 3. Расчет уравнений лучей, определяющих углы обзора виртуальных камер

Выведем уравнения лучей  $l_{ij}$ . Поскольку  $\angle AOB = \angle COB = \alpha$ , то уравнения лучей  $OA$  и  $OC$  имеют вид  $z = x \cdot \cot(\mp\alpha)$ . Точка  $A$  является пересечением луча  $OA$  и прямой, на которой располагается левое зеркало. Координаты точки  $C$  вычисляются по аналогии:

$$(9) \quad A, C = \frac{b}{\cot(\mp\alpha) - \tan\phi_{1,2}} \cdot \left[ \begin{array}{c} 1 \\ \cot(\mp\alpha) \end{array} \right].$$

Уравнения лучей  $l_{ij}$  имеют следующий вид:

$$(10) \quad \begin{aligned} l_{i1} : z &= -\cot(2\phi_{5-i}) \cdot x + b = \tan\left(2\phi_{5-i} + \frac{\pi}{2}\right) \cdot x + b, \\ l_{i2} : z &= \tan\left(2\phi_{5-i} + j\alpha + \frac{\pi}{2}\right) \cdot x + 2b \frac{\cos\phi_{5-i} \sin(\phi_{5-i} + j\alpha)}{\sin(2\phi_{5-i} + j\alpha)}, \\ & j = 7 - 2i. \end{aligned}$$

Для дальнейших выкладок нам потребуется информация об углах наклона прямых  $l_{ij}$ , их можно получить из уравнений прямых, но следует учесть периодичность функции тангенс и набор условий (1). В итоге получаем:

$$(11) \quad \angle l_{ik} = \begin{cases} 2\phi_{5-i} + j\alpha(k-1) + \frac{\pi}{2} - \pi, & \phi_{5-i} \in \left(\frac{\pi}{2}, \frac{3\pi}{4} - j\frac{\alpha}{2}(k-1)\right] \\ 2\phi_{5-i} + j\alpha(k-1) + \frac{\pi}{2} - 2\pi, & \phi_{5-i} \in \left(\frac{3\pi}{4} - j\frac{\alpha}{2}(k-1), \pi\right] \end{cases}.$$

Путем геометрических построений нетрудно убедиться, что первый вариант соответствует случаю, когда луч  $l_{ij}$  направлен вверх от оси  $OX$  и не пересекает ее, а второй вариант — когда направлен вниз и пересекает.

### 4. Ограничения на размер зеркал

Вычислим расстояния  $AB$  и  $CB$  (ширина левого и правого зеркал):

$$(12) \quad AB, CB = \frac{b \sin \alpha}{|\cos(\phi_{1,2} \mp \alpha)|}.$$

С учетом набора условий (1) нетрудно убедиться, что оба выражения под модулями всегда меньше нуля. Пусть максимальная ширина левого и правого зеркал ограничена некоторыми значениями  $L_{10}$  и  $L_{20}$  (например, в силу технологических причин или ограничений на габариты технического средства, на котором установлена оптическая система). В итоге ограничения на размер зеркал задаются следующими неравенствами:

$$(13) \quad \begin{aligned} L_1(\phi_1, b) &:= \frac{b \sin \alpha}{-\cos(\phi_1 - \alpha)} \leq L_{10}, \\ L_2(\phi_2, b) &:= \frac{b \sin \alpha}{-\cos(\phi_2 + \alpha)} \leq L_{20}. \end{aligned}$$

### 5. Условия, при которых не наблюдается взаимное отражение зеркал

Исследуем, при каких условиях в поля зрения виртуальных камер не попадают зеркала. С учетом выражений (1) и (2) получаем, что точка  $T_2$  всегда лежит правее оси  $OZ$ , следовательно, поле обзора правой виртуальной камеры не может пересекать первую координатную четверть, в которой находится правое зеркало.

Перейдем к левой виртуальной камере. Исходя из геометрических соображений, факт наблюдения камерой левого зеркала определяется взаимным положением зеркала и луча  $l_{31}$ . Угол  $\angle l_{31}$  может выражаться двумя разными способами (выражение (11)). Если  $\phi_2 \in [\frac{3\pi}{4}, \pi]$ , то луч  $l_{31}$  будет направлен вниз, а левое зеркало будет располагаться над полем зрения левой виртуальной камеры (в плоскости  $OXZ$ ) и заведомо не будет попадать в это поле. Если же  $\phi_2 \in (\frac{\pi}{2}, \frac{3\pi}{4}]$ , то камера не будет наблюдать левое зеркало при  $\phi_1 < \angle l_{31} = 2\phi_2 - \frac{\pi}{2}$ .

Если же  $\phi_2 \in [\frac{3\pi}{4}, \pi]$ , то  $(2\phi_2 - \frac{\pi}{2}) \in [\pi, \frac{3\pi}{2}]$ , т.е. при таких значениях  $\phi_2$  условие  $(\phi_1 < 2\phi_2 - \frac{\pi}{2})$  тоже выполняется. В итоге, условия, при которых не наблюдается взаимное отражение зеркал, записывается следующим образом:

$$(14) \quad \phi_1 < 2\phi_2 - \frac{\pi}{2}.$$

## 6. Условия, при которых зеркала уместятся внутри короба фиксированного размера

В наших предыдущих статьях [31, 32] исследовался образец двухзеркальной оптической системы — стереонасадка для цифровой камеры Raspberry Pi Camera, данная камера достаточно широко в различных проектах, связанных с компьютерным зрением. Камера предназначена для подключения к микрокомпьютеру Raspberry Pi с помощью шлейфа (рисунок 2). Согласно информации с *официального сайта*<sup>URL</sup>, данная камера (ее первая версия) имеет горизонтальный угол обзора  $53,5^\circ = 0,934$  рад., т.е.  $\alpha = 26,75^\circ = 0,467$  рад.

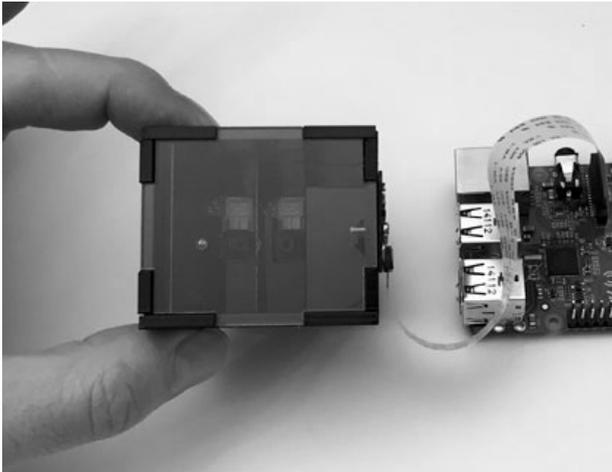


Рисунок 2. Стереонасадка с установленной цифровой камерой Raspberry Pi Camera, подключенной к микрокомпьютеру

Вся конструкция заключена внутри пластикового короба в виде прямоугольного параллелепипеда, одна из граней короба сделана из прозрачного пластика, остальные грани — из черного пластика. Размещение зеркал внутри короба позволяет жестко зафиксировать зеркала относительно камеры и защищает зеркала от внешних воздействий. На рисунке 3 показана схема стереонасадки (вид сверху). Стороны короба — это прямоугольник  $DEWA$ , сторона  $AD$  является прозрачной, а камера крепится на стороне  $DE$ . Короб определяется своей шириной  $w_{box} = \|DE\|$  и высотой  $x_{box} = \|AD\|$ . Центр объектива камеры принимается за начало координат, которое располагается на стороне  $DE$ .

Максимальные размеры подобного короба могут быть ограничены: например, если требуется разместить его на беспилотном летательном

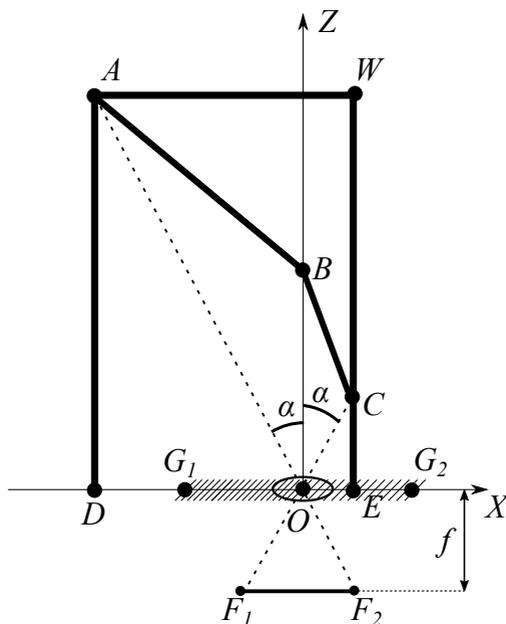


РИСУНОК 3. Схема двухзеркальной системы внутри прямоугольного короба

аппарате или наземном подвижном роботе. Следующие два неравенства обеспечивают то, что конструкция из зеркал будет помещаться внутри короба заданного размера:

$$(15) \quad \begin{aligned} (X_C - X_A) &\leq w_{box}, \\ Z_A &\leq h_{box}. \end{aligned}$$

## 7. Условия, при которых в зеркалах не отражаются запрещенные области

Любая реальная видеокамера имеет определенные габариты, в том числе лицевая часть видеокамеры. При использовании камеры в подобной катоптрической системе возможна ситуация, когда камера будет отражаться в зеркалах, чего желательно избегать, поскольку на изображениях окажутся неинформативные области. Например, камера Raspberry Pi Camera размещена на микросхеме в форме квадрата со стороной 2,5 см., а объектив располагается приблизительно в центре микросхемы. Таким

образом, в поле зрения виртуальных камер не должна попадать не только видеокамера, но и области слева и справа от видеокамеры шириной по 1,25 см. на стороне  $DE$ . Назовем подобные области запрещенными. На рисунке 3 отрезок  $G_1G_2$  соответствует области, которая занимает конструкция видеокамеры.

В общем случае, точки  $G_1$  и  $G_2$  могут располагаться или внутри отрезка  $DE$ , или вне его, или совпадать с точками  $D$  и  $E$ . Если  $G_1$  располагается внутри отрезка  $DO$ , то отрезок  $DG_1$  может быть прозрачным и может попадать в поле зрения виртуальных камер, что продемонстрировано на рисунке 4. Здесь изображена 3D-модель, включающая в себя два зеркала и короб. На грани, включающей в себя отрезок  $G_1E$ , закрепляется видеокамера. Точки  $E$  и  $G_2$  совпадают. То есть, запрещенной областью является не вся сторона  $DE$ , а только ее часть, отрезок  $G_1E$ .

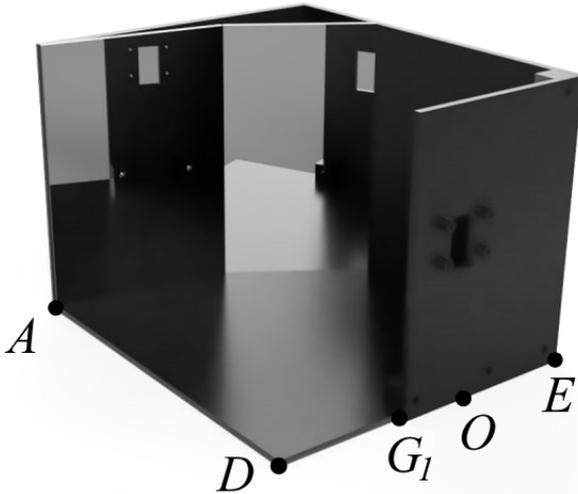


Рисунок 4. 3D-модель катоптрической системы внутри короба с частично прозрачной стороной  $DE$

С учетом вышесказанного, введем новые обозначения: пусть запрещенная область задается в виде отрезка  $S_1O$  длиной  $s_1$  и отрезка  $OS_2$ , оба они лежат на оси  $OX$ . В поле зрения правой виртуальной камеры отрезок  $S_1S_2$  не попадает, если выполняется одно из условий:

- луч  $l_{41}$  не пересекает прямую, на которой лежит этот отрезок. Тогда луч  $l_{41}$  направлен вниз (по отношению к положительному направлению оси  $OZ$ ). Согласно выражению (11), это выполняется,

если угол наклона левого зеркала не превосходит  $\frac{3\pi}{4}$ .

- луч  $l_{41}$  пересекает отрезок  $S_1S_2$  левее точки  $S_1$ ;
- луч  $l_{42}$  пересекает отрезок правее точки  $S_2$  — невозможно при допустимых конфигурациях зеркал.

По аналогичной схеме исследуются условия, при которых отрезок  $S_1S_2$  не попадает в поле зрения левой виртуальной камеры. Обозначим через  $u_{ij}$  абсциссу точки пересечения луча  $l_{ij}$  и оси  $OX$ . В итоге условия, при которых в зеркалах не отражаются запрещенные области, описываются следующим набором выражений:

$$(16) \quad \left( \phi_1 < \frac{3\pi}{4} \right) OR \left( \left( \phi_1 \geq \frac{3\pi}{4} \right) AND (u_{41} < -s_1) \right), \\ \left( \phi_2 < \left( \frac{3\pi}{4} - \frac{\alpha}{2} \right) \right) OR \left( \left( \phi_2 \geq \left( \frac{3\pi}{4} - \frac{\alpha}{2} \right) \right) AND (u_{32} < -s_1) \right).$$

Сторону  $EW$  можно не рассматривать, поскольку ее нижняя часть не попадает в поле зрения видеокамеры, а та часть, которая могла бы попасть, перекрыта зеркалами.

## 8. Постановка и решение задачи условной оптимизации для поиска оптимальной конфигурации оптической системы

Зафиксируем параметры задачи. Будем использовать видеокамеру Raspberry Pi Camera, угол  $\alpha = 26,75^\circ = 0,467$  рад.. Если использовать разрешение  $960 \times 720$  пикс., то фокальная длина  $f$  равна  $960 / (2 \tan \alpha) = 952$  пикс. Обозначим условия, в которых будет функционировать проектируемая оптическая система. Она закреплена на беспилотном летательном аппарате, который будет использоваться для инспекции состояния высотных зданий и линий электропередач. Введем минимально допустимую величину стереобазы, используя выражение (7). Пусть  $Z_0 = 100$  см., примем эту величину в качестве безопасного рабочего расстояния между БПЛА и исследуемыми объектами. Погрешность  $\Delta d$  положим равной 2 пикс. Это значение, конечно же, сильно зависит от особенностей используемого алгоритма поиска стереосоответствия. Допустимую погрешность  $\Delta Z_0$  в вычислении расстояния положим равной 3 см. Тогда  $T_{min} = 7$  см. Максимально возможную ширину обоих зеркал положим равной  $L_{10} = L_{20} = 10$  см. Зеркала заключены внутри корпуса, максимально допустимые размеры которого равны  $wbox = hbox = 15$  см. Камера располагается в центре квадратной микросхемы со стороной 2,5 см. (как уже упоминалось с предыдущем разделе). Пусть ширина  $s_1 = |G_1O|$  запрещенной области равна 1,5 см., немного больше половины

ширины микросхемы, а отрезок  $DG_1$  прозрачен, т.е. может попадать в поле зрения виртуальных камер.

Полный набор ограничений на конфигурацию оптической системы задается выражениями (1), (8), (13), (14), (15), (16). Оптимизируемыми параметрами являются два угла  $\phi_1, \phi_2$  и длина  $b$ , в совокупности они однозначно определяют конфигурацию зеркал относительно видеокамеры. При проектировании оптических систем часто необходимо не только достижение заданных оптических характеристик, но и минимизация размеров системы. Одним из возможных параметров, характеризующих размер рассматриваемой системы, является периметр прямоугольника  $DEWA$ . Т есть, необходимо найти условный минимум следующей функции:

$$(17) \quad F(\phi_1, \phi_2, b) = 2((X_C(\phi_1, b) - X_A(\phi_2, b)) + Z_A(\phi_2, b)) \rightarrow \min.$$

Аналитическое решение задачи получить не удалось, поэтому задача решалась численно, с использованием возможностей программного пакета *SciPy*<sup>URL</sup>. Была задействована функция *differential\_evolution* из модуля *optimize*, она позволяет вычислять глобальный минимум скалярной функции  $F(\vec{X})$  при заданном наборе ограничений вида  $\vec{b}_l \leq \vec{F}(\vec{X}) \leq \vec{b}_u$ , здесь  $\vec{b}_l$  и  $\vec{b}_u$  — векторы из вещественных чисел. Используется алгоритм дифференциальной эволюции, который является эвристическим. При вызове функции значения всех ее настроек были выбраны по умолчанию.

Программа доступна в *github-репозитории*<sup>URL</sup>. В результате был найдено минимальное значение целевой функции, равное 20,6 см., оно достигается при  $\phi_1 = 172^\circ = 3$  рад.,  $\phi_2 = 131^\circ = 2,288$  рад.,  $b = 5,3$  см. Размеры зеркал — 3 см. и 2,6 см. Полученная оптическая система помещается в короб размером  $(X_C - X_A) \times Z_A = 4,6 \times 5,8$  см.

Осталось найти высоту обоих зеркал и короба. Необходимо добиться, чтобы в поле зрения видеокамеры не попадало ничего, кроме зеркал. Будем использовать следующее соображение: точка  $A$  является наиболее удаленной от видеокамеры точкой, принадлежащей зеркалам, поэтому если левый край левого зеркала занимает по высоте весь кадр с видеокамеры, то зеркала будут занимать все поле зрения камеры. Здесь мы также учитываем, что зеркала прямоугольные и параллельны оси  $OY$ . Камера Raspberry Pi Camera имеет вертикальный угол обзора  $2\beta = 41,41^\circ = 0,934$  рад. Тогда высота зеркал равна  $2Z_A \tan(\beta) = 5,9$  см.

Использовать предложенную методику расчета параметров оптической системы можно следующим образом: имеется некоторая видеокамера с известными характеристиками (разрешение и углы обзора). Это соответствует ситуации, когда проектировщик оптической системы не

имеет возможности заниматься разработкой и созданием непосредственно видеокамер, но он может выбрать одну из нескольких доступных видеокамер, чтобы в итоге получить оптическую стереосистему с заданными характеристиками. Заданы максимально допустимые размеры короба, внутри которого должна помещаться оптическая система, и максимально допустимые ширины зеркал. Необходимо рассчитать конфигурацию зеркал таким образом, чтобы минимизировать размеры короба, но при этом чтобы полученное значение стереобазы было не меньше нужного значения. Размеры короба и зеркал могут задаваться, исходя из особенностей конкретной задачи: например, для установке на подвижном роботе вся конструкция должна быть достаточно компактной, а для задачи, описанной, к примеру, в статье [26], ограничения на габариты могут быть менее строгими.

### 9. Расчет коэффициента увеличения для виртуальных камер

Если в зеркалах отражается некоторый точечный объект с координатами  $[X \ Y \ Z]^t$ , то оптический путь до него будет различаться для левой и правой виртуальных камер, т.е., два изображения будут получены с разным коэффициентом увеличения. Оценим, насколько могут различаться эти коэффициенты при различных положениях объекта. Для этого вычислим координаты объекта относительно левой и правой виртуальных камер (выражение (2)) и рассмотрим отношение расстояний от объекта до плоскостей изображений обеих камер:

$$\begin{aligned}
 \check{Z}_{1,2} &= X \sin 2\phi_{2,1} + (b - Z) \cos 2\phi_{2,1} + b, \\
 \tilde{S} &:= \sqrt{X^2 + (b - Z)^2}, \sin \theta := \frac{X}{\tilde{S}}, \cos \theta := \frac{b - Z}{\tilde{S}}, \\
 \hat{S} &:= \frac{\tilde{S}}{b}, \check{Z}_{1,2} = b \left( \hat{S} \cos (2\phi_{2,1} - \theta) + 1 \right), \\
 \frac{\check{Z}_2}{\check{Z}_1} &= K(\hat{S}, \theta) = \frac{K_1(\hat{S}, \theta)}{K_2(\hat{S}, \theta)} = \frac{\hat{S} \cos (2\phi_1 - \theta) + 1}{\hat{S} \cos (2\phi_2 - \theta) + 1} \rightarrow extr, \\
 \theta &\in [0, 2\pi), \hat{S} \geq 0.
 \end{aligned}
 \tag{18}$$

Таким образом, положение объекта задается в полярной системе координат, привязанной к точке соприкосновения зеркал, угол  $\theta$  задается относительно отрицательного направления оси  $OZ$ . В то же время, объект должен находиться в области, образованной пересечением полей зрения обеих виртуальных камер (обозначим эту область как  $V$ ), т.е. есть

ограничения на возможные значения величин  $\widehat{S}$  и  $\theta$ . В программной среде Maple было выполнено исследование возможных конфигураций области  $V$  при различных значениях углов  $\phi_1$  и  $\phi_2$ . Моделирование показало, что область  $V$  может иметь одну из трех конфигураций, в зависимости от значения  $\phi_1 - \phi_2$ :

$$(19) \quad \begin{aligned} a) & (\phi_1 - \phi_2) < \frac{\alpha}{2}, \\ b) & (\phi_1 - \phi_2) \in \left[ \frac{\alpha}{2}, \alpha \right], \\ c) & (\phi_1 - \phi_2) \in \left( \alpha, \frac{\pi}{2} \right). \end{aligned}$$

Результаты расчета оптической системы из предыдущего раздела подпадают под случай  $c$ ). В случае  $a$ ) область  $V$  имеет наиболее простую конфигурацию и лежит между лучами  $l_{31}$  и  $l_{41}$ , в двух остальных случаях конфигурация более сложная и будет подробно рассмотрена в следующей публикации. Оба луча выходят из точки соприкосновения зеркал. Координаты множества точек, лежащих между этими лучами, можно параметризовать следующим образом:  $\widehat{S} \geq 0, \theta \in [\theta_4, \theta_3]$ . Геометрические построения показали, что  $\theta_i = 2\phi_{i-2} - 2\pi$ . Нужно найти экстремальные точки функции  $K(\widehat{S}, \theta)$ :

$$(20) \quad \begin{aligned} (\phi_1 - \phi_2) < \frac{\alpha}{2} \in \left( 0, \frac{\pi}{4} \right) & \Rightarrow \cos 2(\phi_1 - \phi_2) \in (0, 1), \\ \theta \in [2\phi_2 - 2\pi, 2\phi_1 - 2\pi], 2\phi_1 - \theta \in [2\pi, 2(\phi_1 - \phi_2) + 2\pi] & \Rightarrow \\ \cos(2\phi_1 - \theta) > 0 & \Rightarrow K_1(\widehat{S}, \theta) > 0. \end{aligned}$$

Аналогично доказывается, что  $K_2(\widehat{S}, \theta) > 0$ . Таким образом, функция  $K(\widehat{S}, \theta)$  непрерывна и ограничена внутри области  $V$ . Также функция не имеет стационарных точек внутри этой области, следовательно, достигает максимального значения на границах области:

$$(21) \quad \max K(\widehat{S}, \theta) = K(+\infty, \theta_4) = \frac{1}{\cos 2(\phi_1 - \phi_2)}.$$

Рассчитаем коэффициент увеличения для конкретной двухзеркальной системы. Пусть  $\alpha = 40^\circ = 0,698$  рад.,  $\phi_1 = 154^\circ = 2,688$  рад.,  $\phi_2 = 135^\circ = 2,356$  рад.,  $b = 5$  см,  $S_{\max} = 200$  см. Согласно выражению (8), величина стереобазы будет равна 5,2 см,  $K(\widehat{S}_{\max}, \theta_4) = 1,256$ , а  $\max K(\widehat{S}, \theta) =$

1,269.

## Заключение

Разработана математическая модель катоптрической системы, состоящей из двух плоских зеркал и обеспечивающей получение стереопар изображений с помощью единственной видеокамеры. Модель имеет качественные преимущества перед аналогами за счет учета практических аспектов, которые возникают при проектировании подобных систем: ограничения на размер отражающих элементов, возможность взаимного отражения зеркал, возможность отражения в зеркалах запрещенных областей, общие габариты оптической системы.

Выполнена постановка задачи условной оптимизации для поиска оптимальной конфигурации рассматриваемой оптической системы, в качестве целевой функции выбран периметр короба в форме прямоугольного параллелепипеда, в этот короб заключена катоптрическая система, и его периметр необходимо минимизировать. С использованием возможностей программного пакета SciPy была написана программа для численного решения этой задачи. Выполнен теоретический расчет коэффициента увеличения для виртуальных камер.

В дальнейшем будут рассмотрены ограничения, связанные с конфигурацией области, которая является пересечением полей зрения двух виртуальных камер. Полученные ограничения будут добавлены в уже рассмотренную задачу условной оптимизации. Также планируется использовать предложенный подход к исследованию трехзеркальных и четырехзеркальных катоптрических стереосистем.

## Список использованных источников

- [1] Gorevoy A. V., Machikhin A. S. *Optimal calibration of a prism-based videoendoscopic system for precise 3D measurements* // Computer Optics.– 2017.– Vol. 41.– No. 4.– Pp. 535–544.  [↑25](#)
- [2] Zhou F., Chen Y., Zhou M., Li X. *Effect of catadioptric component postposition on lens focal length and imaging surface in a mirror binocular system* // Sensors.– 2019.– Vol. 23.– No. 19.– id. 5309.– 20 pp.  [↑25](#)
- [3] Gluckman J., Nayar S. K. *Rectified catadioptric stereo sensors* // IEEE Transactions on Pattern Analysis and Machine Intelligence.– 2002.– Vol. 24.– No. 2.– Pp. 224–236.  [↑25](#)
- [4] Clark A. F., Chan S. W. *Single-camera computational stereo using a rotating mirror* // Proc. 1994 British Machine Vision Conference, ed. Hancock E. R.– BMVA Press.– 1994.– ISBN 952-1898-1-X.– Pp. 761–770.   [↑25](#)

- [5] Nakao T., Kashitani A. *Panoramic camera using a mirror rotation mechanism and a fast image mosaicing* // *Proc 2001 International Conference on Image Processing* (07–10 October 2001, Thessaloniki, Greece).– 2001.– ISBN 0-7803-6725-1.– Pp. 1045–1048. doi ↑25
- [6] Hu S., Dong H., Shimasaki K., Jiang M., Senoo T., Ishii I. *Omnidirectional panoramic video system with frame-by-frame ultrafast viewpoint control* // *IEEE Robotics and Automation Letters*.– 2022.– Vol. 7.– No. 2.– Pp. 4086–4093. doi ↑25
- [7] Pachidis T., Lygouras J. *A pseudo stereo vision system as a sensor for real time path control of a robot* // *Proceedings of the 19th IEEE Instrumentation and Measurement Technology Conference*.– V. 2, IMTC/2002 (21–23 May 2002, Anchorage, AK, USA).– 2002.– ISBN 0-7803-7218-2.– Pp. 1589–1594. doi ↑25, 26
- [8] Vernon D. *An optical device for computation of binocular stereo disparity with a single static camera*, Opto-Ireland 2002: Optical Metrology, Imaging, and Machine Vision, Proc. SPIE.– vol. 4877.– 2003.– Pp. 38–46. doi ↑25, 26
- [9] Chai X., Zhou F., Chen X. *Epipolar constraint of single camera mirror binocular stereo vision systems* // *Optical Engineering*.– 2017.– Vol. 56.– No. 8.– id. 084103.– 8 pp. doi ↑25, 26
- [10] Zhou F., Chai X., Chen X., Song Y. *Omnidirectional stereo vision sensor based on single camera and catoptric system* // *Applied Optics*.– 2016.– Vol. 55.– No. 25.– Pp. 6813–6820. doi ↑26
- [11] Liu Y., Zhou F., Guo Z., Tan H., Zhang W. *Design and optimization of a quad-directional stereo vision sensor with wide field of view based on single camera* // *Measurement*.– 2022.– Vol. 203.– No. 7.– id. 111915.– 11 pp. doi ↑26
- [12] Wang R., Li X., Zhang Y. *Analysis and optimization of the stereo-system with a four-mirror adapter* // *Journal of the European Optical Society Rapid Publications*.– 2008.– Vol. 3.– id. 08033.– 7 pp. doi ↑26
- [13] Yu L., Pan B. *Structure parameter analysis and uncertainty evaluation for single-camera stereo-digital image correlation with a four-mirror adapter* // *Applied Optics*.– 2016.– Vol. 55.– No. 25.– Pp. 6936–6946. doi ↑26
- [14] Luo H., Yu L., Pan B. *Design and validation of a demand-oriented single-camera stereo-DIC system with a four-mirror adapter* // *Measurement*.– 2021.– Vol. 186.– No. 5.– id. 110083.– 13 pp. doi ↑26
- [15] López-Alba E., Felipe-Sesé L., Schmeer S., Díaz F. A. *Optical low-cost and portable arrangement for full field 3D displacement measurement using a single camera* // *Measurement Science and Technology*.– 2016.– Vol. 27.– No. 11.– id. 115901. doi ↑26
- [16] Yu Z., Ma K., Wang Z., Wu J., Wang T., Zhuge J. *Surface modeling method for aircraft engine blades by using speckle patterns based on the virtual stereo vision system* // *Optics Communications*.– 2018.– Vol. 411.– No. 1.– Pp. 33–39. doi ↑27
- [17] Bartol K., Bojanić D., Petković T., Pribanić T. *Catadioptric stereo on a smartphone* // *Proc. 2021 12th International Symposium on Image and Signal Processing and Analysis*, ISPA (13–15 September 2021, Zagreb, Croatia), Piscataway, NJ: IEEE.– 2021.– ISBN 1-66542-639-X.– Pp. 189–194. doi ↑27
- [18] Aggarwal R., Vohra A., Nambodiri A. M. *Panoramic stereo videos with a single camera* // *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (27–30 June 2016, Las Vegas, NV, USA).– IEEE.– 2016.– ISBN 978-1-4673-8850-4.– Pp. 3755–3763. doi ↑27

- [19] Zhu L., Wang W., Liu Y., Lai S., Li J. *A virtual reality video stitching system based on mirror pyramids // 2017 International Conference on Virtual Reality and Visualization (ICVRV)* (21–22 October 2017, Zhengzhou, China).– IEEE.– 2017.– ISBN 978-1-5386-2636-8.– Pp. 288–293. doi ↑27
- [20] Nene S. A., Nayar S. K. *Stereo with mirrors // Sixth International Conference on Computer Vision* (07 January 1998, Bombay, India).– IEEE.– 1998.– ISBN 81-7319-221-9.– Pp. 1087–1094. doi ↑27
- [21] Baker S., Nayar S. K. *A theory of single-viewpoint catadioptric image formation // International Journal of Computer Vision.*– 1999.– Vol. **35**.– No. 2.– Pp. 175–196. doi ↑27
- [22] Goshasby A., Gruver W. A. *Design of a single-lens stereo camera system // Pattern Recognition.*– 1993.– Vol. **26**.– No. 6.– Pp. 923–937. doi ↑27
- [23] Gluckman J., Nayar S. K. *Planar catadioptric stereo: geometry and calibration // Proceedings 1999 IEEE Computer Society Conference on Computer Vision and Pattern Recognition.*– V. 1 (23–25 June 1999, Fort Collins, CO, USA).– IEEE.– 1999.– ISBN 0-7695-0149-4.– Pp. 22–28. doi ↑27
- [24] Gluckman J., Nayar S. K. *Catadioptric stereo using planar mirrors // International Journal of Computer Vision.*– 2001.– Vol. **44**.– No. 1.– Pp. 65–79. doi ↑27
- [25] Endres F., Sprunk C., Kümmerle R., Burgard W. *A catadioptric extension for RGB-D cameras // 2014 IEEE/RSJ International Conference on Intelligent Robots and Systems* (14–18 September 2014, Chicago, IL, USA).– IEEE.– 2014.– ISBN 9781479969357.– Pp. 466–471. doi ↑27, 28
- [26] Mariottini G. L., Scheggi S., Morbidi F., Prattichizzo D. *Catadioptric stereo with planar mirrors: multiple-view geometry and camera localization // Visual Servoing via Advanced Numerical Methods, Lecture Notes in Control and Information Sciences.*– vol. **401**, eds. Chesi G., Hashimoto K., London: Springer.– 2010.– ISBN 978-1-84996-088-5.– Pp. 3–21. doi ↑28, 43
- [27] Durand-Texte T., Melon M., Simonetto E., Durand S. *3D vision method applied to measure the vibrations of non-flat items with a two-mirror adapter*, 13th International Conference on Vibration Measurements by Laser and Noncontact Techniques (20–22 June 2018, Ancona, Italy) // *Journal of Physics Conference Series.*– 2018.– Vol. **1149**.– No. 1.– id. 012008.– 9 pp. doi ↑28
- [28] Takahashi K., Nobuhara S. *Structure of multiple mirror system from kaleidoscopic projections of single 3D point // IEEE Transactions on Pattern Analysis and Machine Intelligence.*– 2022.– Vol. **44**.– No. 9.– Pp. 5602–5617. doi ↑28
- [29] Zhao Y., Chen Y., Yang L. *Calibration of double-plane-mirror catadioptric camera based on coaxial parallel circles // Journal of Sensors.*– 2022.– Vol. **2022**.– id. 145400.– 15 pp. doi ↑28, 29
- [30] Zhong F, Quan C. *A single color camera stereo vision system // IEEE Sensors Journal.*– 2018.– Vol. **18**.– No. 4.– Pp. 1474–1482. doi ↑28
- [31] Степанов Д. Н., Смирнов А. В. *Исследование процесса калибровки и оптических характеристик стереонасадки 3Dberry // Программные системы: теория и приложения.*– 2018.– Т. **9**.– № 3(38).– С. 11–28. doi URL ↑28, 38
- [32] Степанов Д. Н. *Математические модели получения стереоизображений с двухзеркальных катадиоптрических систем с учетом дисторсии объектов // Компьютерная оптика.*– 2019.– Т. **43**.– № 1.– С. 105–114. doi ↑30, 33, 38

- [33] Bradski G., Kaehler A. *Learning OpenCV: Computer Vision with the OpenCV Library*.— Sebastopol, CA: O'Reilly Media Inc.— 2008.— ISBN 978-0-596-51613-0.— 575 pp. ↑<sup>35</sup>

Поступила в редакцию 18.04.2024;  
 одобрена после рецензирования 17.06.2024;  
 принята к публикации 26.08.2024;  
 опубликована онлайн 10.09.2024.

Рекомендовал к публикации

*д.ф.-м.н. С. М. Абрамов*

### Информация об авторах:



#### Дмитрий Николаевич Степанов

к.т.н., научн. сотр. Исследовательского центра мультипроцессорных систем ИПС им. А.К. Айламазяна РАН. Область научных интересов: математическое моделирование, численные методы, компьютерное зрение, распознавание образов, параллельное программирование, визуальная навигация, анализ данных



0000-0003-2582-5757

*e-mail:*



#### Игорь Петрович Тищенко

к.т.н., и.о. директора ИПС им. А.К. Айламазяна РАН. Область научных интересов: распознавание образов, параллельное программирование, искусственные нейронные сети, беспилотные летательные аппараты



0000-0002-0369-0524

*e-mail:*

Вклад авторов: *Д. Н. Степанов* – 95% (идея, методология, программное обеспечение, валидация, формальный анализ, расследование, сбор материала, курирование данных, написание черновой версии, доработка и редактирование, визуализация); *И. П. Тищенко* – 5% (наставничество, администрирование).

Декларация об отсутствии личной заинтересованности: *благополучие авторов не зависит от результатов исследования.*



# Mathematical modeling and research of the optimal configuration of an optical stereo system consisting of two flat mirrors

Dmitry Nikolaevich **Stepanov**<sup>1</sup>, Igor Petrovich **Tishchenko**<sup>2</sup>

<sup>1,2</sup>Ailamazyan Program Systems Institute of RAS, Ves'kovo, Russia

<sup>1</sup> [mitek1989@mail.ru](mailto:mitek1989@mail.ru)

**Abstract.** The paper is devoted to mathematical modeling and optimization of optical stereo system configuration, consists of video camera and two flat mirrors. The difference between this research and previous researches is the consideration of a large number of restrictions on the configuration of the optical system: the size of the stereo base, the size of the mirrors, overall dimensions of the optical system, the absence of double reflection of light rays, preventing the situation when the video camera is reflected in the mirrors. A conditional optimization problem is formulated to find the optimal configuration of the considered optical system. The perimeter of the rectangle limiting the dimensions of the optical system was chosen as the target function. Numerical solution to the problem was found using the SciPy package. The results obtained expand the theory of computer vision and can be used in the creation and research of computer vision systems for robotic systems. (*In Russian*).

**Key words and phrases:** machine vision, optical devices, mathematical modeling, stereovision, optimization, catoptric system

2020 *Mathematics Subject Classification:* 78M50; 78A05, 68T45

**For citation:** Dmitry N. Stepanov, Igor P. Tishchenko. *Mathematical modeling and research of the optimal configuration of an optical stereo system consisting of two flat mirrors*. Program Systems: Theory and Applications, 2024, **15**:3(62), pp. 23–53. (*In Russ.*). [https://psta.psisiras.ru/read/psta2024\\_3\\_23-53.pdf](https://psta.psisiras.ru/read/psta2024_3_23-53.pdf)

## References

- [1] A. V. Gorevoy, A. S. Machikhin. “Optimal calibration of a prism-based videoendoscopic system for precise 3D measurements”, *Computer Optics*, **41**:4 (2017), pp. 535–544. [doi](#)
- [2] F. Zhou, Y. Chen, M. Zhou, X. Li. “Effect of catadioptric component postposition on lens focal length and imaging surface in a mirror binocular system”, *Sensors*, **23**:19 (2019), id. 5309, 20 pp. [doi](#)
- [3] J. Gluckman, S. K. Nayar. “Rectified catadioptric stereo sensors”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **24**:2 (2002), pp. 224–236. [doi](#)
- [4] A. F. Clark, S. W. Chan. “Single-camera computational stereo using a rotating mirror”, *Proc. 1994 British Machine Vision Conference*, ed. Hancock E. R., BMVA Press, 1994, ISBN 952-1898-1-X, pp. 761–770. [doi](#) [URL](#)
- [5] T. Nakao, A. Kashitani. “Panoramic camera using a mirror rotation mechanism and a fast image mosaicing”, *Proc 2001 International Conference on Image Processing (07–10 October 2001, Thessaloniki, Greece)*, 2001, ISBN 0-7803-6725-1, pp. 1045–1048. [doi](#)
- [6] S. Hu, H. Dong, K. Shimasaki, M. Jiang, T. Senoo, I. Ishii. “Omnidirectional panoramic video system with frame-by-frame ultrafast viewpoint control”, *IEEE Robotics and Automation Letters*, **7**:2 (2022), pp. 4086–4093. [doi](#)
- [7] T. Pachidis, J. Lygouras. “A pseudo stereo vision system as a sensor for real time path control of a robot”, *Proceedings of the 19th IEEE Instrumentation and Measurement Technology Conference. V. 2, IMTC/2002 (21–23 May 2002, Anchorage, AK, USA)*, 2002, ISBN 0-7803-7218-2, pp. 1589–1594. [doi](#)
- [8] D. Vernon. “An optical device for computation of binocular stereo disparity with a single static camera”, *Opto-Ireland 2002: Optical Metrology, Imaging, and Machine Vision*, Proc. SPIE, vol. **4877**, 2003, pp. 38–46. [doi](#)
- [9] X. Chai, F. Zhou, X. Chen. “Epipolar constraint of single camera mirror binocular stereo vision systems”, *Optical Engineering*, **56**:8 (2017), id. 084103, 8 pp. [doi](#)
- [10] F. Zhou, X. Chai, X. Chen, Y. Song. “Omnidirectional stereo vision sensor based on single camera and catoptric system”, *Applied Optics*, **55**:25 (2016), pp. 6813–6820. [doi](#)
- [11] Y. Liu, F. Zhou, Z. Guo, H. Tan, W. Zhang. “Design and optimization of a quad-directional stereo vision sensor with wide field of view based on single camera”, *Measurement*, **203**:7 (2022), id. 111915, 11 pp. [doi](#)
- [12] R. Wang, X. Li, Y. Zhang. “Analysis and optimization of the stereo-system with a four-mirror adapter”, *Journal of the European Optical Society Rapid Publications*, **3** (2008), id. 08033, 7 pp. [doi](#)
- [13] L. Yu, B. Pan. “Structure parameter analysis and uncertainty evaluation for single-camera stereo-digital image correlation with a four-mirror adapter”, *Applied Optics*, **55**:25 (2016), pp. 6936–6946. [doi](#)

- [14] H. Luo, L. Yu, B. Pan. “Design and validation of a demand-oriented single-camera stereo-DIC system with a four-mirror adapter”, *Measurement*, **186**:5 (2021), id. 110083, 13 pp. [doi](#)
- [15] E. López-Alba, L. Felipe-Sesé, S. Schmeer, F. A. Díaz. “Optical low-cost and portable arrangement for full field 3D displacement measurement using a single camera”, *Measurement Science and Technology*, **27**:11 (2016), id. 115901. [doi](#)
- [16] Z. Yu, K. Ma, Z. Wang, J. Wu, T. Wang, J. Zhuge. “Surface modeling method for aircraft engine blades by using speckle patterns based on the virtual stereo vision system”, *Optics Communications*, **411**:1 (2018), pp. 33–39. [doi](#)
- [17] K. Bartol, D. Bojanić, T. Petković, T. Pribanić. “Catadioptric stereo on a smartphone”, *Proc. 2021 12th International Symposium on Image and Signal Processing and Analysis*, ISPA (13–15 September 2021, Zagreb, Croatia), IEEE, Piscataway, NJ, 2021, ISBN 1-66542-639-X, pp. 189–194. [doi](#)
- [18] R. Aggarwal, A. Vohra, A. M. Namboodiri. “Panoramic stereo videos with a single camera”, *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (27–30 June 2016, Las Vegas, NV, USA), IEEE, 2016, ISBN 978-1-4673-8850-4, pp. 3755–3763. [doi](#)
- [19] L. Zhu, W. Wang, Y. Liu, S. Lai, J. Li. “A virtual reality video stitching system based on mirror pyramids”, *2017 International Conference on Virtual Reality and Visualization (ICVRV)* (21–22 October 2017, Zhengzhou, China), IEEE, 2017, ISBN 978-1-5386-2636-8, pp. 288–293. [doi](#)
- [20] S. A. Nene, S.K. Nayar. “Stereo with mirrors”, *Sixth International Conference on Computer Vision* (07 January 1998, Bombay, India), IEEE, 1998, ISBN 81-7319-221-9, pp. 1087–1094. [doi](#)
- [21] S. Baker, S. K. Nayar. “A theory of single-viewpoint catadioptric image formation”, *International Journal of Computer Vision*, **35**:2 (1999), pp. 175–196. [doi](#)
- [22] A. Goshtasby, W. A. Gruver. “Design of a single-lens stereo camera system”, *Pattern Recognition*, **26**:6 (1993), pp. 923–937. [doi](#)
- [23] J. Gluckman, S. K. Nayar. “Planar catadioptric stereo: geometry and calibration”, *Proceedings 1999 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. V. 1 (23–25 June 1999, Fort Collins, CO, USA), IEEE, 1999, ISBN 0-7695-0149-4, pp. 22–28. [doi](#)
- [24] J. Gluckman, S. K. Nayar. “Catadioptric stereo using planar mirrors”, *International Journal of Computer Vision*, **44**:1 (2001), pp. 65–79. [doi](#)
- [25] F. Endres, C. Sprunk, R. Kümmerle, W. Burgard. “A catadioptric extension for RGB-D cameras”, *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems* (14–18 September 2014, Chicago, IL, USA), IEEE, 2014, ISBN 9781479969357, pp. 466–471. [doi](#)
- [26] G. L. Mariottini, S. Scheggi, F. Morbidi, D. Prattichizzo. “Catadioptric stereo with planar mirrors: multiple-view geometry and camera localization”, *Visual Servoing via Advanced Numerical Methods*, Lecture Notes in Control and Information Sciences, vol. **401**, eds. Chesi G., Hashimoto K., Springer, London, 2010, ISBN 978-1-84996-088-5, pp. 3–21. [doi](#)

- [27] T. Durand-Texte, M. Melon, E. Simonetto, S. Durand. “3D vision method applied to measure the vibrations of non-flat items with a two-mirror adapter”, 13th International Conference on Vibration Measurements by Laser and Noncontact Techniques (20–22 June 2018, Ancona, Italy), *Journal of Physics Conference Series*, **1149**:1 (2018), id. 012008, 9 pp. [doi](#)
- [28] K. Takahashi, S. Nobuhara. “Structure of multiple mirror system from kaleidoscopic projections of single 3D point”, *IEEE Transactions on Pattern Analysis and Machine Intelligence*, **44**:9 (2022), pp. 5602–5617. [doi](#)
- [29] Y. Zhao, Y. Chen, L. Yang. “Calibration of double-plane-mirror catadioptric camera based on coaxial parallel circles”, *Journal of Sensors*, **2022** (2022), id. 145400, 15 pp. [doi](#)
- [30] F Zhong, C. Quan. “A single color camera stereo vision system”, *IEEE Sensors Journal*, **18**:4 (2018), pp. 1474–1482. [doi](#)
- [31] D. N. Stepanov, A. V. Smirnov. “Research of calibration process and optical characteristics of 3Dberry stereo nozzle”, *Program Systems: Theory and Applications*, **9**:3(38) (2018), pp. 11–28 (in Russian). [doi](#)
- [32] D. N. Stepanov. “Mathematical models of obtaining stereo images from two-mirror catadioptric systems with regard to lens distortion”, *Computer Optics*, **43**:1 (2019), pp. 105–114 (in Russian). [doi](#)
- [33] G. Bradski, A. Kaehler. *Learning OpenCV: Computer Vision with the OpenCV Library*, O’Reilly Media Inc, Sebastopol, CA, 2008, ISBN 978-0-596-51613-0, 575 pp.

УДК 004.93'11

10.25209/2079-3316-2024-15-3-53-74



# Применение Сиамских нейронных сетей для классификации биомассы растений по визуальному состоянию

Александр Владимирович **Смирнов**<sup>1</sup>, Игорь Петрович **Тищенко**<sup>2</sup>

<sup>1,2</sup>Институт программных систем им. А. К. Айламазяна РАН, Вельсково, Россия

**Аннотация.** В настоящей статье предложен метод классификации биомассы растений по визуальному состоянию с использованием изображений, снятых в специально сконструированной теплице, и технологий искусственных нейронных сетей Сиамской архитектуры. Определены критерии различных состояний биомассы растений. Сформирован собственный набор данных для обучения Сиамских нейронных сетей, содержащий в себе образцы состояний биомассы в форме текстур. В результате была получена точность при обучении в 91.6% и средняя точность классификации отдельных состояний биомассы в 73.6%.

**Ключевые слова и фразы:** Сиамские нейронные сети, набор данных, биомасса растений, классификация

**Для цитирования:** Смирнов А. В., Тищенко И. П. *Применение Сиамских нейронных сетей для классификации биомассы растений по визуальному состоянию* // Программные системы: теория и приложения. 2024. Т. 15. № 3(62). С. 53–74. [https://psta.psir.ru/read/psta2024\\_3\\_53-74.pdf](https://psta.psir.ru/read/psta2024_3_53-74.pdf)

## Введение

На сегодняшний день, в сфере информационных технологий широкое распространение получает использование искусственных нейронных сетей (ИНС). Одним из наиболее популярных направлений применения ИНС является обработка изображений. Обработка и анализ изображений может затрагивать различные сферы жизнедеятельности людей, в том числе и

сельское хозяйство (выращивание различных растений). Применение ИНС для анализа внешнего вида и состояния растений поможет в последствии автоматизировать процесс выращивания растений, что в свою очередь станет неотъемлемой частью автоматических систем для ухода за растениями.

В работе [1] авторы использовали глубинные нейронные сети для оценки динамики роста растений, в частности технологию Social LEAP Estimates Animal Poses (SLEAP)<sup>1</sup>. Оценка происходила по кадрам небольших видеороликов и охватывала различные параметры растений от внешнего вида и формы кроны до условий освещения. В итоге авторы сделали вывод о том, что технология SLEAP достаточно точно отслеживает рост растения в боковой проекции, но имеет проблемы с видом сверху из-за перекрытия листьев.

Работа [2] также посвящена анализу роста растений. В данной работе авторы анализировали рост различных видов салата и использовали несколько архитектур свёрточных нейронных сетей (СНС) таких, как: FFNN<sup>2</sup> и ViLSTM<sup>3</sup>. Исследование предполагало периодический сбор изображений растений в течении 4-х недельного периода. Всего для обучения нейронных сетей было собрано 443 изображения трех сортов салата. В результате, авторам удалось добиться точности по R2<sup>4</sup> до 80%. Авторы отмечают, что на итоговую точность повлияло малое количество изображений.

Кроме анализа роста растений, СНС используется и для их классификации. Например, в работе [3] используется нейросетевую модель CNN-DFLC, которая предназначена для анализа структуры листьев растений и выполнения классификации. В данном исследовании для обучения СНС был использован набор данных VNPlant-200 [4], состоящий из 20000 изображений в 200 классах. Таким образом, авторами была получена средняя точность классификации равная 96,42%.

Другой вид анализа растений предполагает нахождение и классификацию их болезней. Работа [5] представляет собой обширный обзор различных методов определения и классификации болезней растений, в том числе с использованием СНС. В данном обзоре уделяется внимание болезням, поражающим 11 различных растений, которые можно детектировать визуально. Приведено описание и классификация болезней растений и факторов их возникновения. В результате своего исследования, авторы

---

<sup>1</sup>*Social LEAP Estimates Animal Poses (SLEAP)*<sup>URL</sup>

<sup>2</sup>*Feedforward neural network*<sup>URL</sup>

<sup>3</sup>*Bidirectional LSTM*<sup>URL</sup>

<sup>4</sup>*Оценка R2 в машинном обучении*<sup>URL</sup>

сформировали ряд ограничивающих факторов для использования СНС для определения и классификации болезней растений, среди которых особенно выделяются проблемы с количеством и качеством обучающих данных/изображений.

В статье [6] описывается применение глубинных СНС для распознавания болезней листьев растений. В данном исследовании были испытаны 18 различных нейросетевых архитектур. В итоге, авторы предлагают собственную нейросетевую модель PlaNet, точность которой сравнима с аналогами и составляет 96.86%.

Аналогичная работа представлена в [7]. Здесь авторы также используют нейронные сети глубокого обучения для классификации болезней листьев растений. В этой статье исследуется новая модель СНС с несколькими автоматическими экстракторами признаков, а именно СНС плотного слияния (DFNet). В качестве средства извлечения признаков использовались модели MobileNetV2<sup>5</sup> и NASNetMobile<sup>6</sup>. Предложенный метод был протестирован на наборах данных болезней листьев кукурузы и кофе, где была достигнута точность в 97.53% и 94.65% соответственно.

Также технологии нейронных сетей применяются для оценки биомассы растений. Например, в статье [8] рассматривается применение двух методов машинного обучения (многомерная регрессионная сеть<sup>7</sup> и нейронная сеть на основе ResNet-50<sup>8</sup>) для прогнозирования роста биомассы растений. В данной работе авторы используют набор данных, сформированный из изображений 57 растений, снятых с двух разных ракурсов в течение пяти дней. В итоге наилучшие оценки биомассы были получены с помощью многомерной регрессионной сети, что дало среднеквадратичную ошибку 0.0466. Наилучшие оценки относительной скорости роста были получены с помощью сети ResNet-50, что дало среднеквадратичную ошибку 0.1767.

Помимо использования классических СНС, для анализа изображений растений можно использовать Сиамские нейронные сети. Сиамские нейронные сети – это особый класс нейронных сетей, основной задачей которого является не непосредственная классификация образов, а выявление сходства или различия между входными данными. К особенностям Сиамских нейронных сетей можно отнести использование относительно небольших объёмов данных для обучения.

Среди более ранних (2019 – 2020 г.) работ можно выделить [9], [10] и [11]. В данных статьях Сиамские нейронные сети используются для классификации листьев растений, их видов и болезней. Во всех перечисленных

---

<sup>5</sup>MobileNet, MobileNetV2, and MobileNetV3<sup>URL</sup>

<sup>6</sup>NasNetLarge and NasNetMobile<sup>URL</sup>

<sup>7</sup>Deep Learning Models for Multi-Output Regression<sup>URL</sup>

<sup>8</sup>ResNet (34, 50, 101): «остаточные» CNN для классификации изображений<sup>URL</sup>

работах была достигнута точность классификации сравнимая с точностью при использовании глубинных нейронных сетей.

Подобные исследования также представлены в статье [12], где Сиамская нейронная сеть используется вместе с СНС для распознавания видов растений и выявления болезней. Авторы создают гибридную модель (нейросетевой ансамбль), в котором Сиамская сеть отвечает за определения вида растения, а СНС модель VGG16<sup>9</sup> за распознавание заболелания. Такой подход позволил увеличить точность распознавания на 13.4%–37.39% по сравнению с использованием одиночных моделей СНС.

Если рассматривать Сиамские нейронные сети в качестве инструмента для анализа изображений наряду с глубинными СНС, то можно обратиться к работе [13]. Здесь Авторы проводят исследование с целью доказать, что использование сиамских сетей может быть эффективнее с точки зрения затрат времени без потери точности классификации. В качестве подсетей использовалась сеть LeNet-5<sup>10</sup>. Обучение проводилась на наборах данных MNIST<sup>11</sup>, Fashion-MNIST<sup>12</sup> и CIFAR10<sup>13</sup>. Точность составила 99.11%, 91.65% и 81.64%, что сопоставимо с точностью при классическом использовании сети LeNet-5 при двукратном уменьшении времени обучения.

Исходя из вышеизложенного, можно сделать вывод о том, что в данный момент активно ведутся работы по применению искусственных нейронных сетей (как свёрточных, так и Сиамских) для анализа изображений растений, что подтверждает актуальность настоящей работы.

Настоящая статья посвящена разработке метода классификации биомассы растений по визуальному состоянию с использованием Сиамских нейронных сетей. Будет представлено определение визуального состояния биомассы, а также описан метод формирования обучающего набора данных. Проведено экспериментальное тестирование разработанного метода, в результате которого была достигнута средняя точность в 73.6% по метрике F-score.

## 1. Постановка задачи

При проектировании систем автоматизированного ухода за растениями очень важно уделить внимание сенсорным подсистемам, которые необходимым для сбора информации об окружающей среде и о состоянии растений. Однако, однозначное определение состояния растения в каком-либо

---

<sup>9</sup>VGG16 — нейросеть для выделения признаков изображений<sup>URL</sup>

<sup>10</sup>Архитектура LeNet-5 с использованием Python<sup>URL</sup>

<sup>11</sup>THE MNIST DATABASE<sup>URL</sup>

<sup>12</sup>Fashion MNIST<sup>URL</sup>

<sup>13</sup>The CIFAR-10 dataset<sup>URL</sup>

числовом эквиваленте – задача нетривиальная и не может быть решена посредством использования датчика/сенсора, как, например, определение влажности почвы или температуры воздуха.

Выходом из ситуации мог бы стать анализ внешнего вида или визуального состояния растений. Чтобы его выполнить, необходимо получить качественное изображение исследуемого растения. Однако, установка множества цифровых камер в теплицу для наблюдения за растениями финансово не выгодно, так как несёт затраты на покупку самих камер, и обеспечения их влагозащиты и электропитания. К тому же, могут появиться проблемы связанные с недостатком освещения или неверным ракурсом съёмки.

Тем не менее состояние растения – это тот показатель, который может свидетельствовать об эффективности автоматизированной системы ухода и по сути является некоторым видом «обратной связи» растения с системой. В таком случае, следует использовать не анализ внешнего вида отдельных растений, а анализ визуального состояния биомассы в целом. Для этого достаточно одной камеры установленной в верхней части теплицы, объектив которой будет направлен вертикально вниз.

Таким образом, задача настоящего исследования заключается в сегментации/классификации запечатлённой на снимке биомассы растений на классы, различающиеся по состояниям растения. В проведённом исследовании рассматривается только два состояния биомассы растения: «здоровое растение» и «больное растение». Следовательно, необходимо определить к какому классу относится входное изображение биомассы растения, что можно сделать, сравнив его с образцом одного из классов. Для решения подобных задач применяется специальная архитектура нейронных сетей – Сиамские нейронные сети.

Целью же настоящего исследования является получение работоспособного метода для классификации биомассы растений по визуальному состоянию с использованием технологий Сиамских нейронных сетей, для последующего применения в автоматизированных системах ухода и выращивания растений, в качестве сенсорной подсистемы. Также следует отметить тот факт, что проводимое исследование не было направленно на классификацию отдельных видов растений или распознавание их заболеваний.

## **2. Критерии состояния растений**

В рамках настоящего исследования рассматриваются только два крайних состояния биомассы растений: «здоровое растение» и «больное растение». Сокращение возможных состояний до двух крайних, было

вызвано тем, что достаточно затруднительно однозначно дифференцировать промежуточные состояния растения на используемых изображениях (рисунок 1).



Рисунок 1. Пример используемых изображений

Критерии определения состояния биомассы были сформированы исходя из внешнего вида наблюдаемых растений. Стоит отметить тот факт, что в настоящем исследовании наблюдались растения с визуально схожим жизненным циклом, такие как: *Огурец*, *Томат*, *Баклажан*, *Перец сладкий*, *Перец острый*. Данные растения, находясь в оптимальном состоянии, имеют ровные листья зелёного или оттеков зелёного цвета, но при недостаточном питании или наличии иных вредоносных факторов цвет листьев изменяется в сторону жёлто-коричневых оттенков, на листьях появляются повреждения, которые можно определить визуально. Таким образом, в основном оценивался цвет и форма листьев наблюдаемых растений. Если они имели свой естественный цвет и форму без видимых изменений, то состояние биомассы растения определялось как «здоровое». Однако, если по какому-то из критериев растение не проходило, то оно считалось «больным». На рисунке 2 показан биомассы «здорового» и «больного» растения.

Определение состояния биомассы растения и соответствия его внешнего вида сформированным ранее критериям происходило в ручном режиме. Оператор с использованием специального программного обеспечения выделял участки входных изображений, которые, по его мнению, содержали биомассу «здорового» или «больного» растения. В итоге был создан набор



(а) изображение «здорового» растения (б) изображение «больного» растения

Рисунок 2. Пример «здорового» и «больного» растения

изображений, содержащий примеры биомассы «здоровых» и «больных» растений. Более подробно о создании используемого набора данных будет описано в п. 4.1 настоящей статьи.

### 3. Архитектура и особенности Сиамских нейронных сетей

Нейронные сети Сиамской архитектуры были впервые представлены в начале 1990-х годов в работе [14] для решения проблемы проверки подлинности подписей как задачи сопоставления изображений. Сиамские нейронные сети (Siamese Neural Network, SNN) — это класс архитектур нейронных сетей, предназначенных для сравнения и измерения сходства между парами входных выборок. Термин «сиамский» происходит от идеи, что архитектура сети состоит из парных нейронных сетей (часто свёрточных), которые идентичны по структуре и имеют одинаковый набор весовых коэффициентов. Каждая сеть обрабатывает одну входную выборку из пары, а их выходные данные сравниваются для определения сходства или различия между двумя экземплярами входных данных.

Сиамские сети предназначены для решения задач, где прямое обучение с помеченными выборками ограничено или затруднено, поскольку сеть после обучения способна различать похожие и непохожие экземпляры, не требуя явных меток классов.

Архитектура сиамской сети обычно состоит из трех основных компонентов: используемая подсеть, метрика сходства и функция ошибки (рисунок 3).

Подсеть является основным компонентом архитектуры Сиамской сети. Она отвечает за извлечение ключевых признаков и особенностей

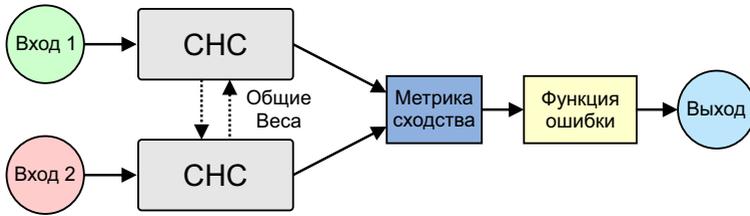


Рисунок 3. Пример архитектуры Сиамской нейронной сети

из входных выборок. Обычно подсети представляют собой свёрточные нейронные сети, состоящие из слоёв свёртки и/или полносвязных слоёв, которые обрабатывают входные данные и создают некоторый дескриптор исследуемого объекта. Распределяя одинаковые веса между идентичными подсетями, модель учится извлекать схожие признаки и особенности для аналогичных входных данных, что позволяет эффективно их сравнивать.

Метрика сходства используется для сравнения сгенерированных дескрипторов и измерения сходства или различия между двумя входными данными. Выбор метрики сходства зависит от конкретной задачи и характера входных данных. Обычно в качестве метрики сходства используется евклидово расстояние, косинусное сходство или коэффициент корреляции.

В роли функции ошибки применяется функция контрастных потерь. Это функция, основанная на подсчёте расстояния, в отличие от более традиционных функций прогнозирования ошибки. Данная функция используется для анализа входных данных, при котором две схожие точки имеют малое евклидово расстояние, а две различные точки имеют большое евклидово расстояние.

Функция контрастных потерь высчитывается по формуле:

$$(1 - Y) \frac{1}{2} (D_W)^2 + (Y) \frac{1}{2} \{\max(0, m - D_W)\}^2$$

где  $Y$  – метка соответствия классов (0 – один класс, 1 – разные),  $m$  – значение предела,  $D_W$  (евклидово расстояние), которое определяется по формуле:

$$\sqrt{\{G_W(X_1) - G_W(X_2)\}^2}$$

где  $X_1$  и  $X_2$  входные данные (изображения),  $G_W$  – выход сети.

Среди преимуществ Сиамских нейронных сетей можно выделить следующее:

- ✓ Нет необходимости в большом наборе данных для обучения. Сеть способна обучиться на небольшом наборе данных, что также позволяет компенсировать дисбаланс классов.

- ✓ Устойчивость к аффинным преобразованиям изображений, таким как поворот и масштабирование.
- ✓ Семантическое сходство. Нейронная сеть Сиамской архитектуры анализирует пространство признаков, чтобы сформировать представление о схожести/различии изображений вместо того, чтобы просто извлекать статические признаки с помощью операции свёртки.

К недостаткам Сиамских нейронных сетей можно отнести:

- Требуется больше времени на обучение по сравнению с традиционными свёрточными нейронными сетями.
- Не предоставляет данные о вероятности определения объекта к какому-либо классу.

#### 4. Реализация и обучение нейронной сети Сиамской архитектуры

Для выполнения поставленных задач по классификации биомассы растений по визуальному состоянию была использована нейронная сеть Сиамской архитектуры с двумя идентичными свёрточными подсетями, которые содержат три слоя свёртки, три полносвязных слоя и имеют следующую конфигурацию:

```
(cnn1): Sequential
Conv2d(3, 96, kernel_size=(11, 11), stride=(4, 4))
ReLU(inplace=True)
MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1,
  ceil_mode=False)
Conv2d(96, 256, kernel_size=(5, 5), stride=(1, 1))
ReLU(inplace=True)
MaxPool2d(kernel_size=2, stride=2, padding=0, dilation=1,
  ceil_mode=False)
Conv2d(256, 384, kernel_size=(3, 3), stride=(1, 1))
ReLU(inplace=True)
(fc1): Sequential
Linear(in_features=384, out_features=1024, bias=True)
ReLU(inplace=True)
Linear(in_features=1024, out_features=256, bias=True)
ReLU(inplace=True)
Linear(in_features=256, out_features=2, bias=True)
```

Программирование архитектуры используемой Сиамской нейронной сети, её обучение и тестирование происходили с использованием языка программирования Python и фреймворка PyTorch<sup>14</sup>. PyTorch — это

---

<sup>14</sup>PyTorch *GET STARTED*<sup>URL</sup>

фреймворк предназначенный для машинного обучения. Он включает в себя набор инструментов для работы с моделями, используется в обработке естественного языка, компьютерном зрении и других похожих направлениях.

#### 4.1. Создание обучающего набора данных

В свободном доступе существуют наборы данных, содержащих в себе как изображения «здоровых» растений, так и растений, поражённых каким-либо заболеванием. Одним из таких наборов данных является New Plant Diseases Dataset<sup>15</sup>, который содержит в себе 87000 изображений в 38 классах. При проведении предварительных экспериментальных исследований, на этом наборе данных была обучена СНС. Однако, в процессе классификации отдельных растений, каких-либо адекватных результатов достичь не удалось. Вероятно, на результат повлияли различные условия съёмки. На рисунке 4 показан пример изображения *Перца* из набора данных New Plant Diseases Dataset и на используемых снимках.



РИСУНОК 4. Пример изображений листа *Перца* из набора данных New Plant Diseases Dataset и на используемых снимках

Невооружённым глазом видно отличие представленных данных. Изображение из набора более высокого качества и содержит одиночный лист растения. Тогда как на снимках из теплицы запечатлена биомасса растения, и само изображение более низкого качества. В связи с этим, было принято решение о создании собственного небольшого набора данных, содержащего образцы состояний биомассы растений.

<sup>15</sup>[New Plant Diseases Dataset](https://www.github.com/eurobot/new-plant-diseases-dataset)<sup>url</sup>

Определение состояния биомассы растения в соответствии со сформированными ранее критериям происходило в ручном режиме, с использованием специально разработанного программного обеспечения. Данное ПО<sup>16</sup> имеет графический интерфейс пользователя, и позволяет выбрать интересующую область изображения с помощью полигонального выделения, задать название класса и при необходимости отзеркалить полученные образцы. На рисунке 5 изображён интерфейс используемого ПО.



РИСУНОК 5. Интерфейс ПО Marker Image: 1 – меню для загрузки изображения; 2 – загруженное изображение; 3 – кнопка выбора директории для сохранения образцов; 4 – поле для введения названия класса; 5 – выбор метода отзеркаливания образцов; 6 – кнопка подтверждающая разметку выбранной области; 7 – поле показа координат точек полигонального выделения размеченной области; 8 – выбор цвета разметки

Созданный набор данных состоял из образцов размером 100x100 пикселей, распределённых по трём классам:

- «plant» – фрагменты биомассы «здоровых» растений. 211 образцов;
- «dise» – фрагменты биомассы «больных» или «умерших» растений. 219 образцов;
- «back» – элементы фона. 222 образца.

На рисунке 6 показаны примеры образцов каждого из классов.

Предложенный набор данных не содержит изображения конкретных видов растений или их заболеваний, а скорее отражает внешний вид

<sup>16</sup>Свидетельство № 2023680292 о государственной регистрации программы «Marker Image v.1.0». Зарегистрировано в Реестре программ для ЭВМ. Дата регистрации 28.09.2023<sup>URL</sup>

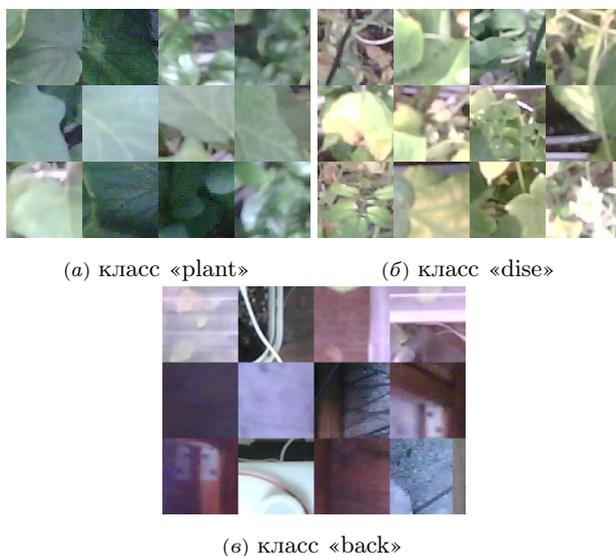


Рисунок 6. Примеры образцов созданного набора данных

текстуры<sup>17</sup> биомассы «здоровых» и «больных» растений, а также фона.

Несмотря на то, что Сиамские сети способны адекватно обучаться при дисбалансе классов, было принято решение о балансировании классов, что потенциально может увеличить точность обучения сети. Процедура балансирования заключалась в добавлении недостающих экземпляров класса, которыми являлись зеркальные копии случайных уже существующих экземпляров. Итого в каждом классе количество экземпляров было выровнено до 230: 200 на обучение и 30 на тест.

#### 4.2. Обучение Сиамской нейронной сети

Обучение используемой Сиамской нейронной сети происходило со следующими параметрами:

- `batch_size`: 16 – размер (количество) данных, посылаемых на вход сети каждую эпоху;
- `epochs`: 100 – количество эпох обучения нейронной сети;
- функция потерь: `ContrastiveLoss` – функция контрастных потерь;
- оптимизатор: `Adam` – один из методов оптимизации обучения, включённых в фреймворк.

<sup>17</sup>[Image texture<sup>URL</sup>](https://www.image-texture.com/)

На рисунке 7 приведён график изменения ошибки обучения в зависимости от эпох.

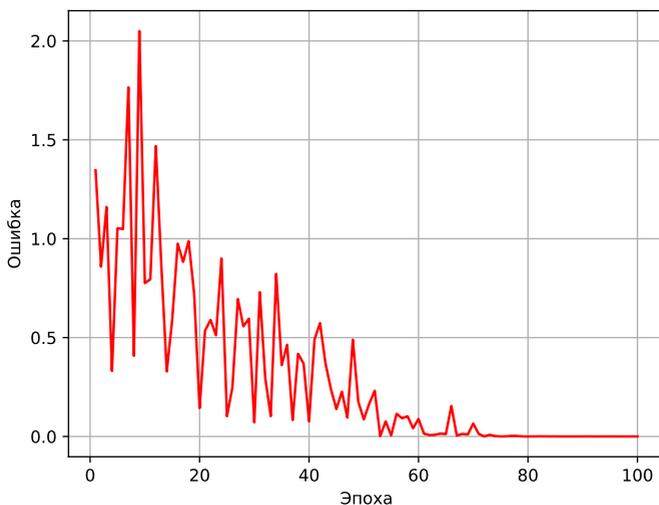


Рисунок 7. График изменения ошибки обучения

### 4.3. Подсчёт точности обучения

В качестве метрики подсчёта точности обучения Сиамской нейронной сети был использован показатель F-score<sup>18</sup>, который высчитывается по формуле:

$$(1) \quad Fscore = \frac{2 * Recall * Precision}{Recall + Precision}$$

$$(2) \quad Precision = \frac{True\ Positive}{True\ Positive + False\ Positive}$$

$$(3) \quad Recall = \frac{True\ Positive}{True\ Positive + False\ Negative}$$

(4)

Здесь следует ввести следующие обозначения:

True Positive (TP) – истинно-положительное решение. Искомый объект обнаружен.

True Negative (TN) – истинно-отрицательное решение. Объект, который не является искомым не был обнаружен.

<sup>18</sup>Что такое F-score и для чего он используется?<sup>url</sup>

False Positive(FP) – ложно-положительное решение. Объект, который не является искомым был детектирован как искомый.

False Negative(FN) – ложно-отрицательное решение. Объект, который является искомым не был обнаружен.

Precision (точность) – отношение TP к TP + FP. Это доля объектов, названными классификатором положительными и при этом действительно являющимися положительными.

Recall (полнота) – отношение TP к TP + FN. Это то, какую долю объектов положительного класса из всех объектов положительного класса нашёл алгоритм.

Поскольку Сиамская нейронная сеть на выходе даёт информацию об относительном расстоянии объектов друг от друга, что является в некотором понимании мерой схожести (меньше расстояние – более похожие объекты), то для определения принадлежности объектов к одному классу было введено пороговое значение расстояния. Если расстояние между объектами ниже данного порога, то они считаются принадлежащими одному классу, в противном случае разным классам. На графике (рисунок 8) отображены пороговые значения с шагом 0.1 и полученная точность обучения, рассчитанная на тестовой выборке.

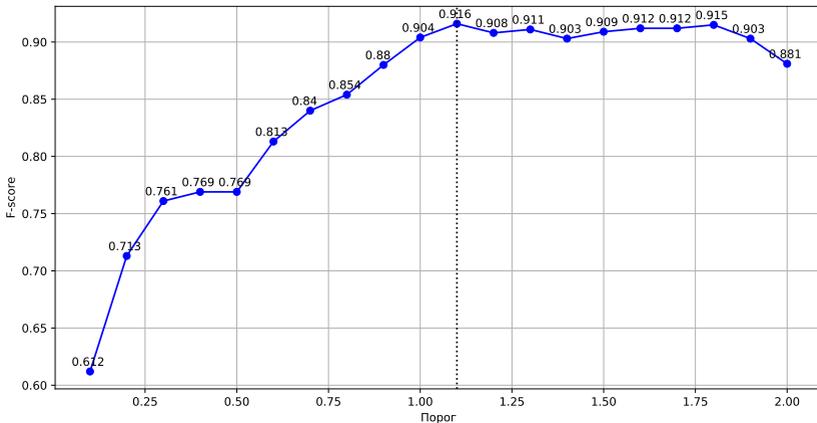


Рисунок 8. График изменения точности от значения порога

## 5. Тестирование обученной Сиамской нейронной сети

В процессе тестирования предложенного метода классификации биомассы растений по визуальному состоянию использовались несколько тестовых изображений, не используемых в качестве обучающих данных.

Для каждого тестового изображения были созданы эталонные данные. Более подробно о создании эталонных данных описано в п. 4.1 настоящей статьи.

### 5.1. Создание эталонных данных

При создании эталонных данных были использованы 5 изображений биомассы растений, которые не задействовались для создания данных для обучающихся. Эти изображения были получены при тех же условиях, что и обучающие данные, то есть сняты с использованием программируемой Wi-Fi камера в специально сконструированной теплице.

Эталонные данные представляли собой бинарные изображения-маски, где интересующая область была окрашена в белый цвет, а всё остальное изображение в чёрный. Для каждого снимка были сгенерированы по 3 изображения-маски в соответствии с количеством исследуемых классов («plant», «dise» и «back»). Генерация изображений-масок происходила в ручном режиме с использованием графического редактора с возможностью редактирования слоёв (например, онлайн редактор Photopea<sup>19</sup>). На рисунке 9 представлен пример снимка биомассы растений и созданных изображений-масок для каждого из классов.

Таким образом, были получены 15 изображений-масок, по 3 маски для каждого из входных изображений биомассы растений. Именно это изображения-маски были использованы в качестве эталонных данных для подсчёта точности классификации биомассы растений по визуальному состоянию.

### 5.2. Классификации растений по их визуальному состоянию

Экспериментальное тестирование классификации биомассы растений по визуальному состоянию происходило на тех же снимках, что использовались для создания эталонных данных. При тестировании был применён подход на основе сканирующего окна<sup>20</sup>. Размер сканирующего окна был получен экспериментально и составил 50 на 50 пикселей. Шаг окна также был получен экспериментально и составлял 1/5 от его размера по горизонтали и вертикали.

В процессе тестирования, область исходного изображения под сканирующим окном подавалась на вход обученной Сиамской нейронной сети в паре с эталонным экземпляром класса, выбранным случайно из тестовой выборки. Затем сеть рассчитывала показатель «непохожести» (dissimilarity)

---

<sup>19</sup>Photopea<sup>URL</sup>

<sup>20</sup>Принцип сканирующего окна<sup>URL</sup>



(a) исходное изображение



(б) изображение-маска для класса «plant» (в) изображение-маска для класса «dise» (г) изображение-маска для класса «back»

Рисунок 9. Исходное изображение биомассы растений и изображения-маски для каждого из классов

для пары входных данных, который показывал на сколько образец схож с эталоном. Если показатель *dissimilarity* был ниже некоторого порогового значения (подбиралось экспериментально), то образец считался экземпляром того же класса, что и эталон, и его координаты на исходном изображении, а также его ширина и высота сохранялись в отдельный список для последующего подсчёта точности классификации.

Для удобства восприятия, область под классифицируемым образцом окрашивалась в определённый цвет (у каждого класса был свой цвет), тем самым создавая некую маску, которая в свою очередь содержала результирующую область класса и накладывалась на копию входного изображения. Таким образом, были сгенерированы изображения (рисунок 10), наглядно показывающие зоны, принадлежащие каждому из классов.



(а) область класса «plant»



(б) область класса «dise»



(в) область класса «back»

РИСУНОК 10. Пример изображений с результирующими областями после классификации

### 5.3. Подсчёт точности классификации растений

Точность классификации биомассы растений рассчитывалась с использованием метрики F-score, которая была описана в п. 4.3 настоящей статьи. В данном случае объект (ранее классифицируемая область изображения с координатами и размером) считался верно определённым (True Positive) если площадь пересечения с эталоном занимала более 50% его собственной площади. В противном случае он относился к необнаруженным объектам (False Negative). Если пересечение между текущим объектом и эталоном отсутствовало, засчитывалось ложно-положительное решение (False Positive), то есть объект, который не является искомым был детектирован как искомый. В таблице 1 представлены результаты подсчёта точности классификации.

ТАБЛИЦА 1. Точность классификации по метрике F-score

Класс	Среднее, %	Худшее, %	Лучшее, %
plant	65	48	71
dise	82	76	91
back	74	68	81

Экспериментальное тестирование обученной нейронной сети показало среднюю точность в 65%, 82% и 74% для классов «plant», «dise» и «back» соответственно. Общая точность по F-score составила около 73.6%.

Несмотря на относительно низкую общую точность, при классификации биомассы «больных растений» было достигнуто значение точности в 91%. Определение наличия «больных растений» и площади их биомассы относительно всей биомассы растений на изображении, по крайней мере, может использоваться как некий детектор, сообщающий о возникновении негативных факторов при выращивании растений. В свою очередь, для нахождения площади всей биомассы растений можно воспользоваться методом, описанным в статье [15]. Тем самым можно будет рассчитать соотношение «здоровых» и «больных растений» к общей биомассе, и на основании полученных значений принять меры по урегулированию негативных влияний.

Также следует отметить тот факт, что на итоговую точность классификации могли повлиять эталонные данные, которые были созданы вручную и не имеют 100% достоверности из-за проблемы однозначной дифференциации растений на используемых изображениях.

### Вывод

В результате проведённого экспериментального исследования был разработан метод классификации биомассы растений по визуальному состоянию с применением нейронных сетей Сямской архитектуры. Была получена точность при обучении равная 91.6% и точность определения биомассы «больных растений» от 76% до 91%.

Представленный метод показал относительно низкое значение точности в сравнении с более классическим применением нейронных сетей. Однако, здесь стоит учитывать различие поставленных задач, условий проведения исследований и сложности используемых нейронных сетей.

Тем не менее, разработанный метод может быть применён в качестве одной из сенсорных подсистем, используемых в системах автоматизированного ухода за растениями, и, как минимум, сигнализировать о наличии больных растений, что позволит своевременно менять правила ухода и корректировать воздействие на растения.

### Список использованных источников

- [1] Gall G. E. C., Pereira T. D., Jordan A., Meroz Y. *Fast estimation of plant growth dynamics using deep neural networks* // Plant Methods.– 2022.– Vol. **18**.– id. 21.– 11 pp. [doi](#) ↑<sup>54</sup>
- [2] Taewon M., Woo-Joo C., Se-Hun J., Da-Seul C., Myung-Min O. *Growth analysis of plant factory-grown lettuce by deep neural networks based on automated feature extraction* // Horticulturae.– 2022.– Vol. **8**.– No. 12.– id. 1124.– 9 pp. [doi](#) ↑<sup>54</sup>
- [3] Savitha P., Mungamuri S. *Accurate plant species analysis for plant classification using convolutional neural network architecture* // International Journal of Reconfigurable and Embedded Systems (IJRES).– 2024.– Vol. **13**.– No. 1.– Pp. 160–170. [doi](#) ↑<sup>54</sup>
- [4] Quoc T. N., Hoang V. T. *VNPlant-200 — A public and large-scale of Vietnamese medicinal plant images dataset* // *Integrated Science in Digital Age 2020, ICIS 2020, Lecture Notes in Networks and Systems*.– vol. **136**, Cham: Springer.– 2020.– ISBN 978-3-030-49263-2.– Pp. 406–411. [doi](#) ↑<sup>54</sup>
- [5] Joseph D. S., Pawar P. M., Pramanik R. *Intelligent plant disease diagnosis using convolutional neural network: a review* // *Multimedia Tools and Applications*.– 2023.– Vol. **82**.– No. 14.– Pp. 21415–21481. [doi](#) ↑<sup>54</sup>
- [6] Khanna M., Singh L. K., Thawkar S., Goyal M. *PlaNet: a robust deep convolutional neural network model for plant leaves disease recognition* // *Multimedia Tools and Applications*.– 2024.– Vol. **83**.– No. 2.– Pp. 4465–4517. [doi](#) ↑<sup>55</sup>
- [7] Faisal M., Leu J. S., Avian C., Prakosa S. W., Köppen M. *DFNet: Dense fusion convolution neural network for plant leaf disease classification* // *Agronomy Journal*.– 2024.– Vol. **116**.– No. 3.– Pp. 826–838. [doi](#) ↑<sup>55</sup>
- [8] Åström O., Hedlund H., Sopsakis A. *Machine-learning approach to non-destructive biomass and relative growth rate estimation in aeroponic cultivation* // *Agriculture*.– 2023.– Vol. **13**.– No. 4.– id. 801.– 13 pp. [doi](#) ↑<sup>55</sup>
- [9] Wang B., Wang D. *Plant leaves classification: a few-shot learning method based on Siamese network* // *IEEE Access*.– 2019.– Vol. **7**.– Pp. 151754–151763. [doi](#) ↑<sup>55</sup>
- [10] Figueroa-Mata G., Mata-Montero E. *Using a convolutional Siamese network for image-based plant species identification with small datasets* // *Biomimetics*.– 2020.– Vol. **5**.– No. 1.– id. 8.– 17 pp. [doi](#) ↑<sup>55</sup>
- [11] Goncharov P., Uzhinskiy A., Ososkov G., Nechaevskiy A., Zudikhina J. *Deep Siamese networks for plant disease detection*, *Mathematical Modeling and Computational Physics 2019 (MMCSP 2019)* // *EPJ Web Conf*.– 2020.– Vol. **226**.– id. 03010.– 4 pp. [doi](#) ↑<sup>55</sup>

- [12] Sherly K. K., Sonia A. *Hybrid CNN models for plant species recognition and disease detection*, Intelligent Computing, Lecture Notes in Networks and Systems.– vol. **1016**, Cham: Springer.– 2024.– ISBN 978-3-031-62280-9.– Pp. 35–50.   ↑56
- [13] Du J., Fu W., Zhang Y., Wang Z. *Advancements in image recognition: a Siamese network approach* // Information Dynamics and Applications.– 2024.– Vol. **3**.– No. 2.– Pp. 89–103.   ↑56
- [14] Bromley J., Guyon I., LeCun Y., Säckinger E., Shah R. *Signature verification using a "Siamese" time delay neural network* // International Journal of Pattern Recognition and Artificial Intelligence.– Vol. **07**.– No. 04.– Pp. 669–688.   ↑59
- [15] Смирнов А. В., Иванов Е. С. *Анализ изображений растения, полученные с камеры системы автоматизированного ухода, для визуальной оценки изменения его состояния с течением времени* // Программные системы: теория и приложения.– 2023.– Т. **14**.– № 3(58).– С. 37–58.     ↑70

Поступила в редакцию 24.06.2024;  
 одобрена после рецензирования 11.08.2024;  
 принята к публикации 11.08.2024;  
 опубликована онлайн 10.09.2024.

Рекомендовал к публикации

*д.ф.-м.н. А. М. Елизаров*

## Информация об авторах:



**Александр Владимирович Смирнов**

Младший научный сотрудник Лаборатории методов обработки и анализа изображений, Институт Программных Систем имени А. К. Айламазяна РАН. Научные интересы: компьютерное зрение; нейронные сети; робототехника; автоматизация и управление

 0000-0002-7104-1462

e-mail:



**Игорь Петрович Тищенко**

Кандидат технических наук, зав. Лабораторией методов обработки и анализа изображений, Институт Программных Систем имени А. К. Айламазяна РАН. Научные интересы: компьютерное зрение; нейронные сети; робототехника; автоматизация и управление

 0000-0002-0369-0524

e-mail:

Вклад авторов: *А. В. Смирнов* – 90% (идея, методология, программное обеспечение, валидация, формальный анализ, расследование, сбор материала, курирование данных, написание черновой версии, доработка и редактирование); *И. П. Тищенко* – 10% (наставничество, администрирование, финансирование).

Декларация об отсутствии личной заинтересованности: *благополучие авторов не зависит от результатов исследования.*



# Application of Siamese neural networks to classify plant biomass by visual state

Alexander Vladimirovich **Smirnov**<sup>1</sup>, Igor Petrovich **Tishchenko**<sup>2</sup>

<sup>1,2</sup>Ailamazyan Program Systems Institute of RAS, Ves'kovo, Russia

**Abstract.** This paper proposes a method for classifying plant biomass by visual condition using images captured in a specially designed greenhouse and Siamese architecture artificial neural network technologies. Criteria for various states of plant biomass have been determined. We have generated our own dataset for training Siamese neural networks, containing samples of biomass states in the form of textures. As a result, a training accuracy of 91.6% and an average classification accuracy of individual biomass states of 73.6%. (*In Russian*).

**Key words and phrases:** Siamese neural networks, dataset, plant biomass, classification

2020 *Mathematics Subject Classification:* 68T10; 68T45,68T07

**For citation:** Alexander V. Smirnov, Igor P. Tishchenko. *Application of Siamese neural networks to classify plant biomass by visual state*. Program Systems: Theory and Applications, 2024, **15**:3(62), pp. 53–74. (*In Russ.*).  
[https://psta.psiras.ru/read/psta2024\\_3\\_53-74.pdf](https://psta.psiras.ru/read/psta2024_3_53-74.pdf)

## References

- [1] G. E. C. Gall, T. D. Pereira, A. Jordan, Y. Meroz. “Fast estimation of plant growth dynamics using deep neural networks”, *Plant Methods*, **18** (2022), id. 21, 11 pp. 
- [2] M. Taewon, C. Woo-Joo, J. Se-Hun, C. Da-Seul, O. Myung-Min. “Growth analysis of plant factory-grown lettuce by deep neural networks based on automated feature extraction”, *Horticulturae*, **8**:12 (2022), id. 1124, 9 pp. 
- [3] P. Savitha, S. Mungamuri. “Accurate plant species analysis for plant classification using convolutional neural network architecture”, *International Journal of Reconfigurable and Embedded Systems (IJRES)*, **13**:1 (2024), pp. 160–170. 

- [4] T. N. Quoc, V. T. Hoang. “VNPlant-200 — A public and large-scale of Vietnamese medicinal plant images dataset”, *Integrated Science in Digital Age 2020*, ICIS 2020, Lecture Notes in Networks and Systems, vol. **136**, Springer, Cham, 2020, ISBN 978-3-030-49263-2, pp. 406–411. 
- [5] D. S. Joseph, P. M. Pawar, R. Pramanik. “Intelligent plant disease diagnosis using convolutional neural network: a review”, *Multimedia Tools and Applications*, **82**:14 (2023), pp. 21415–21481. 
- [6] M. Khanna, L. K. Singh, S. Thawkar, M. Goyal. “PlaNet: a robust deep convolutional neural network model for plant leaves disease recognition”, *Multimedia Tools and Applications*, **83**:2 (2024), pp. 4465–4517. 
- [7] M. Faisal, J. S. Leu, C. Avian, S. W. Prakosa, M. Köppen. “DFNet: Dense fusion convolution neural network for plant leaf disease classification”, *Agronomy Journal*, **116**:3 (2024), pp. 826–838. 
- [8] O. Åström, H. Hedlund, A. Sopasakis. “Machine-learning approach to non-destructive biomass and relative growth rate estimation in aeroponic cultivation”, *Agriculture*, **13**:4 (2023), id. 801, 13 pp. 
- [9] B. Wang, D. Wang. “Plant leaves classification: a few-shot learning method based on Siamese network”, *IEEE Access*, **7** (2019), pp. 151754–151763. 
- [10] G. Figueroa-Mata, E. Mata-Montero. “Using a convolutional Siamese network for image-based plant species identification with small datasets”, *Biomimetics*, **5**:1 (2020), id. 8, 17 pp. 
- [11] P. Goncharov, A. Uzhinskiy, G. Ososkov, A. Nechaevskiy, J. Zudikhina. “Deep Siamese networks for plant disease detection”, Mathematical Modeling and Computational Physics 2019 (MMCP 2019), *EPJ Web Conf.*, **226** (2020), id. 03010, 4 pp. 
- [12] K. K. Sherly, A. Sonia. “Hybrid CNN models for plant species recognition and disease detection”, Intelligent Computing, Lecture Notes in Networks and Systems, vol. **1016**, Springer, Cham, 2024, ISBN 978-3-031-62280-9, pp. 35–50. 
- [13] J. Du, W. Fu, Y. Zhang, Z. Wang. “Advancements in image recognition: a Siamese network approach”, *Information Dynamics and Applications*, **3**:2 (2024), pp. 89–103. 
- [14] J. Bromley, I. Guyon, Y. LeCun, E. Säckinger, R. Shah. “Signature verification using a “Siamese” time delay neural network”, *International Journal of Pattern Recognition and Artificial Intelligence*, **07**:04, pp. 669–688. 
- [15] A. V. Smirnov, E. S. Ivanov. “Analysis of images of a plant obtained from a camera of an automated care system to visually assess the change in its state over time”, *Program Systems: Theory and Applications*, **14**:3(58) (2023), pp. 37–58 (In Russian).  



# Recovering text sequences using deep learning models

Igor Victorovich **Vinokurov**<sup>1</sup>

<sup>1</sup> Financial University under the Government of the Russian Federation, Moscow, Russia

<sup>1</sup> [igvvinokurov@fa.ru](mailto:igvvinokurov@fa.ru)

**Abstract.** This article presents the results of the formation, training and performance evaluation of models with the Encoder-Decoder and Sequence-To-Sequence (Seq2Seq) architectures for solving the problem of supplementing incomplete texts. Problems of this type often arise when restoring the contents of documents from their low-quality images. The studies conducted in the work are aimed at solving the practical problem of forming electronic copies of scanned documents of the «Roskadastr» PLC, the recognition of which is difficult or impossible with standard means.

The formation and study of models was carried out in Python using the high-level API of the Keras package. A dataset consisting of several thousand pairs was formed for the purpose of training and studying the models. Each pair in this set represented an incomplete and corresponding full text. To evaluate the quality of the models, the values of the loss function and the accuracy, BLEU and ROUGE-L metrics were calculated. Loss and accuracy made it possible to evaluate the effectiveness of the models at the level of predicting individual words. The BLEU and ROUGE-L metrics were used to evaluate the similarity between the full and reconstructed texts. The results showed that both the Encoder-Decoder and Seq2Seq models cope with the task of reconstructing text sequences from their fixed set, but the Seq2Seq transformer-based model achieves better results in terms of training speed and quality. (*Linked article texts in English and in Russian*).

**Key words and phrases:** deep learning models, encoder-decoder, sequence-to-sequence transformer, text recovering, BLEU, ROUGE-L, Keras, Python

2020 *Mathematics Subject Classification:* 68T20; 68T07, 68T45

**For citation:** Igor V. Vinokurov. *Recovering text sequences using deep learning models*. Program Systems: Theory and Applications, 2024, **15**:3(62), pp. 75–110. (*In English, in Russian*). [https://psta.psisaras.ru/read/psta2024\\_3\\_75-110.pdf](https://psta.psisaras.ru/read/psta2024_3_75-110.pdf)

## Introduction

In recent years, deep learning models (*Deep Neural Network*, DNN) have achieved significant results in the field of natural language processing (*Neural Language Processing*, NLP) [1]. Analysis of literature sources showed that the most common models used in such tasks as text transformation (translation), text recovery from distorted or incomprehensible documents, scanned documents of poor quality, illegible manuscripts, blurry or damaged images, etc. are Encoder-Decoder and Seq2Seq transformers.

The Encoder-Decoder architecture, based on recurrent neural networks (*Recurrent Neural Networks*, RNN) or convolutional neural networks (*Convolutional Neural Network*, CNN), consists of two main components – encoder and decoder [2]. The encoder transforms the input data into an internal representation taking into account key features of its content. The decoder uses this representation to generate output data by sequentially predicting its elements.

In contrast, the Seq2Seq transformer architecture [3] offers an alternative approach to representing sequences. It uses a transformation mechanism based on multiple layers with attention mechanisms [4], which allows this model to efficiently process long text sequences. The attention mechanism allows the model to take into account the importance of individual words in the context of the entire sentence, thereby contributing to the generation of higher-quality and coherent text. Compared with Encoder-Decoder, the Seq2Seq transformer has several advantages, such as better ability to handle long sequences, a more flexible architecture, and the ability to train models on large amounts of data.

To evaluate the quality of NLP models, one can use both conventional metrics loss, accuracy etc., and metrics specific to evaluating the quality of the generated text, the main ones being BLEU (*Bilingual Evaluation Understudy*) [5] and ROUGE-L (*Recall-Oriented Understudy for Gisting Evaluation – Longest Common Subsequence*) [6]. The first of them measures the similarity between the predicted and reference text. It uses syntactic information to compare sequences of  $n$  words ( $n$ -grams). The more matches in  $n$ -grams between the predicted and reference texts, the higher the BLEU

Сведения об уточняемых земельных участках и их частях						
Сведения о характерных точках границы уточняемого земельного участка с кадастровым номером XXX-XX-XXXX-XXX						
Обозначения характерных точек границы	Существующие координаты, м		Уточненные координаты, м		Средняя квадратическая погрешность положения характерных точек границы, мм	Описание закрепления точки
1	2	3	4	5	6	7

FIGURE 1. Document with highlighted sections of text. Text fragments recognized by OCR are highlighted in color

value. However, this metric does not take into account the semantic and contextual relationship between words, which may limit its applicability. The second ROUGE-L metric evaluates the quality of automatic text summarization. It compares the length of the longest common word sequence between the predicted and reference text with the length of the reference text, thereby measuring the coverage of the predicted text relative to the reference text and allows one to estimate the degree of information compression in the generated text.

The basis for conducting the research, the results of which are presented in this article, was the impossibility of restoring text on scanned documents of poor quality using modern OCR systems, Figure 1.

An obvious solution to this problem is to develop a simple sentence matching system. However, as noted above, DNN models are able to learn complex non-linear dependencies between input and output data, which allows them to more effectively model the context and semantics of text. In addition, DNN models can be more flexible and generalize, which makes them more effective when working with different types of text and information recovery tasks. *It is the generalizing properties of DNN models that served as the rationale for their use in solving the stated problem – a restored and semantically close sentence is better than its complete absence.*

Section 1 provides a rationale for the need for research and a statement of the problem. Section 2 is devoted to the analysis of works on restoring

text sequences. Features of creating a dataset for training models and researching their work are described in section 3. The creation and research of Encoder-Decoder and Seq2Seq models are given in section 4 and section 5 respectively. The advantages and disadvantages of restoring text sequences using the created models are given in section 6.

## 1. Settings the goal and research problems

When recognizing images of text documents using standard OCR tools, it is not always possible to obtain a high-quality copy in text format. The reason is blurry, illegible, or noisy (for example, in the form of handwritten notes) sections of text on the document image [7, 8].

*The goal of this work* is to research the applicability of the main DNN models for NLP – Encoder-Decoder and Seq2Seq transformer for restoring text from a given dataset.

*The goal set in the work can be achieved by solving the following main problems:*

- (1) Creation of a dataset, consisting of preparing pairs in the form of incomplete and corresponding full text.
- (2) Creation of optimal Encoder-Decoder and Seq2Seq transformer models for achieving the goal.
- (3) Analysis of the quality of the models as a result of calculating the loss function and the accuracy, BLEU and ROUGE-L metrics.

## 2. Analysis of works on restoration of text sequences

An analysis of available publications has shown that there are two main types of models that can be used to match and restore text sequences.

- (1) Attention-based models (in most cases, these are Seq2Seq transformer models) that can focus on the context of a sentence and select the most appropriate option for replacing or restoring missing words.
- (2) RNN- and CNN-based models that take into account the context of a sentence and, in the case of CNN, the features of its visualization on images.

Below is a description of the most significant, in the author's opinion, works on the use of these types of models for restoring the original text.

In [9], the authors consider and demonstrate the effectiveness of two classes of attention mechanisms for matching texts in different languages. The first takes into account all source words, the second considers only a subset of the source words.

Missing word recovery using models based on variational autoencoders (VAE) is proposed, for example, in [10]. Here, VAEs implement prediction of the input text sequence for subsequent improvement of RNN training.

The Seq2Seq model, which allows to perform the task of generating corrected text using the attention mechanism, is described in [11]. The model can be used to transform an incorrect, incorrect sequence of characters into corrected text.

A model based on the Seq2Seq architecture with an attention mechanism for error correction in text obtained from an OCR system is given in [13]. The model is trained on a sufficiently large number of text pairs recognized using OCR.

A study of the applicability of BERT (Bidirectional Encoder Representations from Transformers) models for error correction in sentences of non-English native speakers is given in [14]. The model is able to use the context of a sentence to correct grammatical and stylistic errors in the text.

A pre-trained model with the architecture of a masked Seq2Seq transformer for reconstructing a text fragment from its remaining part is considered in [15]. The model's encoder receives a sentence with a randomly masked fragment (several consecutive tokens) as input, and its decoder tries to predict this masked fragment. The paper shows that as a result of fine-tuning, the model is able to reconstruct the original text quite accurately.

The authors of the paper [16] propose a method for correcting errors in text using a neural network based on symbolic self-attention. The model uses the character level to more accurately correct typos, spelling, and other types of errors in the text.

In [17], the author describes the use of RNN for language models and proposes methods for generating text in a given context. The publication is devoted to controlled sequence labelling – an important area of machine learning, including such tasks as speech recognition, handwriting recognition, and part-of-speech labelling.

The study of the applicability of neural network models for restoring distorted character images is carried out in [18]. The authors describe methods for restoring damaged text based on a combination of CNN and RNN.

The article [19] presents a CNN-based model for automatically correcting typos in text. The model is trained on a large corpus of texts and is able to correct typos with high accuracy.

In [20], a method for restoring distorted images using CNN without the need for training on a large dataset is presented. The method was originally developed for image restoration, but can also be applied to restoring damaged text.

The authors of the article [21] present an autoencoder architecture for error correction in OCR systems. The autoencoder model is used to extract hidden text representations and then restore them.

In addition to the above, a fairly large number of works on restoring text sequences are devoted to restoring texts on different types of images of historical documents damaged over time. Their detailed analysis is given in [22].

### 3. Dataset creation

To train the Encoder-Decoder and Seq2Seq models and to research their operation, a proprietary dataset was created, which, as noted above, is a set of pairs of the form:

incomplete text  $\rightarrow$  [start]full text[end]

The main types of documents of the «Roscadastre» PLC software and software package were used to form the dataset. All unique sentences of these documents in the created dataset represent the full text. The

Полный текст (электронный документ)  
Сведения об уточняемых земельных участках и их частях

Неполный текст (результат распознавания OCR)

1. Сведения об
2. об уточняемых
3. Сведения об уточняемых
4. Сведения об уточняемых земельных
5. об уточняемых земельных
6. Сведения об уточняемых земельных участках
7. уточняемых земельных участках
8. Сведения об уточняемых земельных участках и их
9. земельных участках и их частях
10. об уточняемых земельных участках и

...

FIGURE 2. Full text and its corresponding first 10 incomplete texts

results of removing continuous sequences of 1 to  $N - 3$  words from the full text form a set of corresponding incomplete texts (here  $N$  is the number of words in the text). The removal of continuous sequential words from the text is explained by the specifics of blurring and the location of illegible sections of text on scanned documents, see Figure 1.

The number of full-text sentences included in the created dataset does not exceed 250, the number of corresponding incomplete sentences included in the dataset is about 3000. An example of the correspondence between them is shown in Figure 2.

Tokenization and vectorization of text pairs was carried out using `TextVectorization` from the `Keras` package. To remove ambiguity when restoring full texts, text standardization before tokenization involved the inclusion of a feature of the document area in which it may be present.

For training the DNN, 80% of the pairs from the created dataset were used, for validation and testing – 10%.

#### 4. Creation and research of the Encoder-Decoder model

The Encoder-Decoder model was developed and studied in Python using the API `Keras` [23–25]. Figure 3 shows the optimal structure of the model obtained as a result of experimental studies.

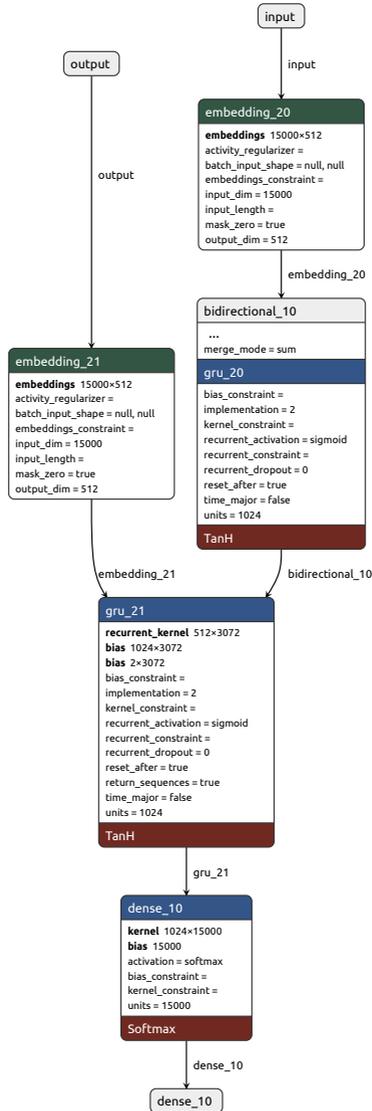


FIGURE 3. Encoder-Decoder model in Keras API for partial text recovery

The encoder of this model consists of the following layers:

- (1) Input layer (`InputLayer`), takes input data in the form of text sequences.
- (2) Embedding layer (`Embedding`) transforms words of text sequences into their vector representations.
- (3) Recurrent layer (`BidirectionalGRU`) processes sequences in both forward and backward directions and produces a hidden state (context vector).

Decoder model layers:

- (1) Input layer (`InputLayer`) takes input as a sequence of words. Determines the shape and type of the input data.
- (2) The embedding layer (`Embedding`) transforms the input tokens into dense vector representations of a given dimensionality, similar to the encoder.
- (3) The recurrent layer (`GRU`) takes the vector representations and processes them sequentially, generating the full text, taking into account the context from the encoder. It is initialized with the state generated by the encoder.
- (4) The output layer (`Dense`) takes the output from the decoder at the current time step and transforms it into a vector of probabilities, each component of which corresponds to a possible next token of the reconstructed text.

The description of the main parameters of all layers of this model is given in table 1.

To assess the quality of the model at the level of individual words of text sequences, the values of the loss function and the accuracy metric were calculated. Figure 4 and Figure 5 show the values of the loss function `sparse_categorical_crossentropy` and the metric `accuracy` of the model for 30 training epochs. The number of epochs was found experimentally and is optimal. The stochastic optimization algorithm `rmsprop` was used to train the model.

The average values of the BLEU and ROUGE-L metrics, which allow us to evaluate the accuracy of reconstructing incomplete text from the testing dataset, are shown in Figure 6. Values in the range of 0.3-0.4 correspond to understanding and acceptable translation of the text. To calculate the metric values, the functions `sentence_bleu()` and `get_scores()` from the NLTK and Rouge packages were used, respectively. They were called after the completion of each training epoch from the callback functions of the `fit()` method.

TABLE 1. Layers of the Encoder-Decoder Model

Layer type (title in Figure 3)	Activity function	Input tensor	Output tensor
<b>InputLayer</b> (input)	–	[(None, None)]	[(None, None)]
<b>Embedding</b> (embedding_20)	–	(None, None)	(None, None, 512)
<b>Bidirectional(GRU)</b> (bidirectional_10)	tanh	(None, None, 512)	(None, 1024)
<b>InputLayer</b> (output)	–	[(None, None)]	[(None, None)]
<b>Embedding</b> (embedding_21)	–	(None, None)	(None, None, 512)
<b>GRU</b> (gru_21)	tanh	[(None, None, 512), (None, 1024)]	(None, None, 1024)
<b>Dense</b> (dense_10)	softmax	(None, None, 1024)	(None, None, 15000)

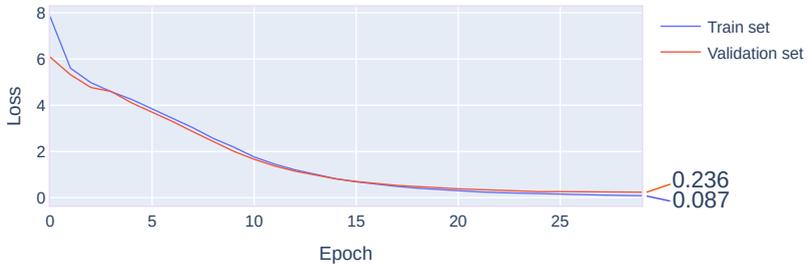


FIGURE 4. Model Losses

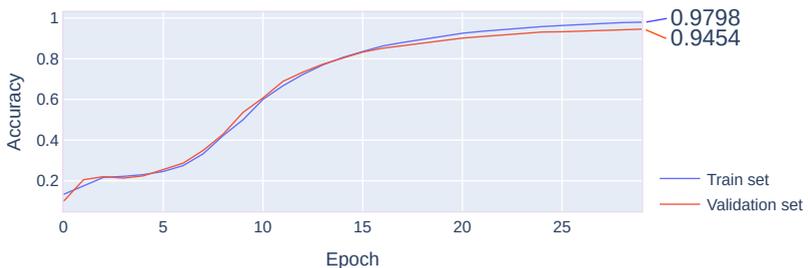


FIGURE 5. Model accuracy

Figure 7 shows the results of text recovery from 10 randomly generated datasets for testing. Each test set consisted of 250–300 sentences. The

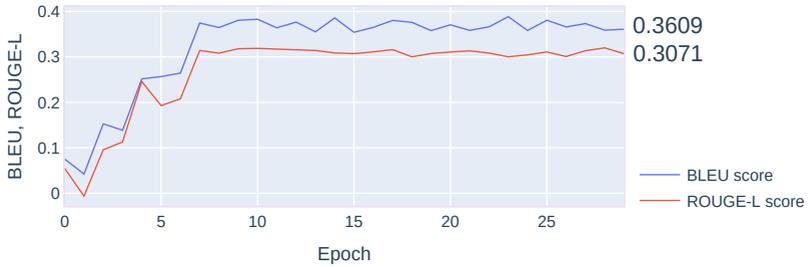


FIGURE 6. Metrics for evaluating the quality of text recovery

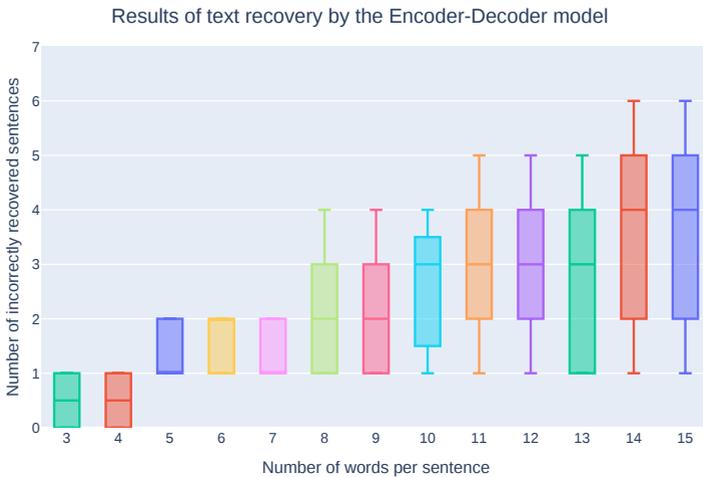


FIGURE 7. Numbers of incorrectly recovered sentences by the Encoder-Decoder model from 10 test datasets

number of incorrectly reconstructed sentences consisting of 11-15 words ranged from 1 to 6.

## 5. Creation and research of the Seq2Seq transformer model

As with the previous model, the creation and research of the model with the Seq2Seq architecture was carried out in Python using the API Keras [23–25].

The optimal structure of the Seq2Seq model, obtained as a result of experimental studies, is shown in Figure 8.

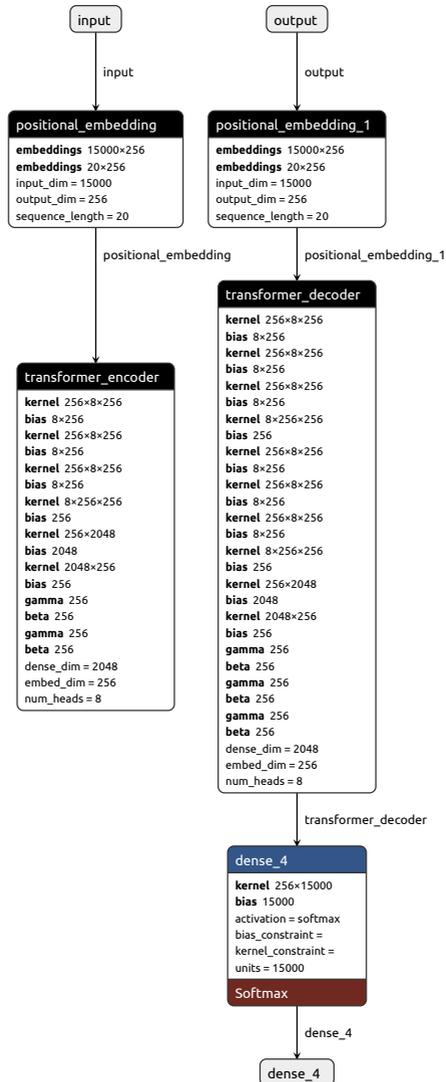


FIGURE 8. Encoder and decoder of the Seq2Seq transformer model in the Keras API for partial text recovery

TABLE 2. Layers of the Seq2Seq transformer model

Layer type (title in Figure 8)	Activity function	Input tensor	Output tensor
<b>InputLayer</b> (input)	–	[(None, None)]	[(None, None)]
<b>PositionalEmbedding</b> (positional_embedding)	–	(None, None)	(None, None, 256)
<b>TransformerEncoder</b> (transformer_encoder)	relu	(None, None, 256)	(None, None, 256)
<b>InputLayer</b> (output)	–	[(None, None)]	[(None, None)]
<b>PositionalEmbedding</b> (positional_embedding_1)	–	(None, None)	(None, None, 256)
<b>TransformerDecoder</b> (transformer_decoder)	relu	(None, None, 256)	(None, None, 256)
<b>Dense</b> (dense_4)	softmax	(None, None, 256)	(None, None, 15000)

The transformer model consists of the following layers:

- (1) The Input Layer (**InputLayer**) takes input as a sequence of words. It determines the shape and type of the input data.
- (2) Positional Embedding Layer (**Positional Embedding Layer**) adds information about the position of a word in a sequence. This allows the model to take into account the order of words in the input and output sequences.
- (3) Transformer Encoder Layers (**TransformerEncoder Layer**), each of which implements an attention mechanism and contains fully connected layers. Due to this, they allow modeling dependencies in sequences, extracting features from the input data and representing them in an optimal internal representation for more complex computations and natural language processing tasks.
- (4) Transformer Decoder Layers (**TransformerDecoder Layer**). Similar to the encoder, the decoder consists of several transformer decoder layers, which also include an attention mechanism and fully connected layers. The decoder generates an output sequence based on the context representations of the encoder.
- (5) The output layer (**Dense**) transforms the predicted token representations into a probability distribution over all possible tokens.

To create more efficient and related representations of text sequences, the attention mechanism [4] is implemented in the encoder and decoder layers of the Seq2Seq architecture model. Table 2 provides a description of the main parameters of all layers of this model.

The study of the accuracy of the model, by analogy with the previous

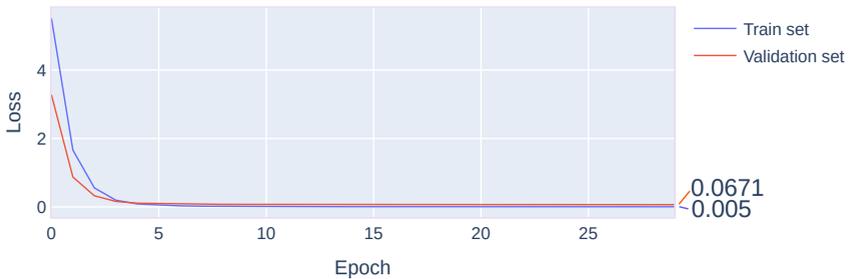


FIGURE 9. Model Losses

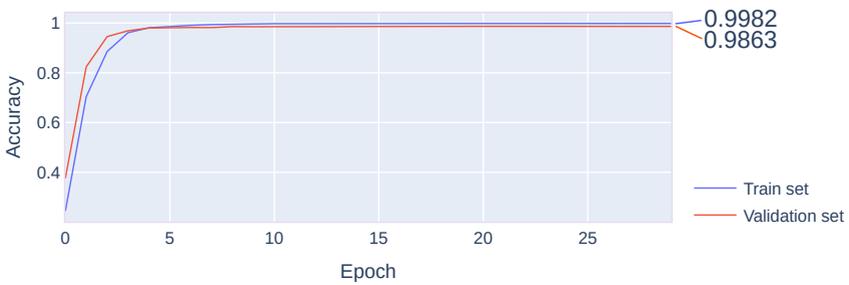


FIGURE 10. Model accuracy

one, consisted in calculating the values of the loss function and the accuracy metric – `sparse_categorical_crossentropy` and `accuracy` respectively for each of the 30 epochs of its training, Figure 9, 10. The number of epochs was found experimentally and is optimal. When training the model, the stochastic optimization algorithm `rmsprop` was used.

The average values of the BLEU and ROUGE-L metrics, which allow us to evaluate the accuracy of the model in reconstructing the incomplete text from the testing dataset, are shown in Figure 11. Values from the range of 0.5-0.6 correspond to high quality of text translation.

Figure 12 shows the results of incomplete text recovery from 10 randomly generated testing datasets. The number of incorrectly recovered sentences is less than for the previous model and ranges from 1 to 4. As a result, it is worth noting the quite natural result – the Seq2Seq transformer model is more efficient than the Encoder-Decoder model due to the use of the attention mechanism. The attention mechanism allows the transformer to highlight key elements of input and output sequences and more effectively identify long-term dependencies in them, which makes the model capable of learning on complex text generation (in this case, text recovery) tasks.

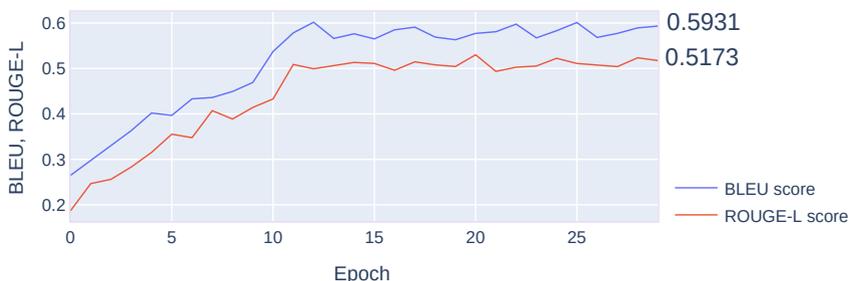


FIGURE 11. Metrics for evaluating the quality of text recovery

Results of text recovery by the Seq2Seq transformer

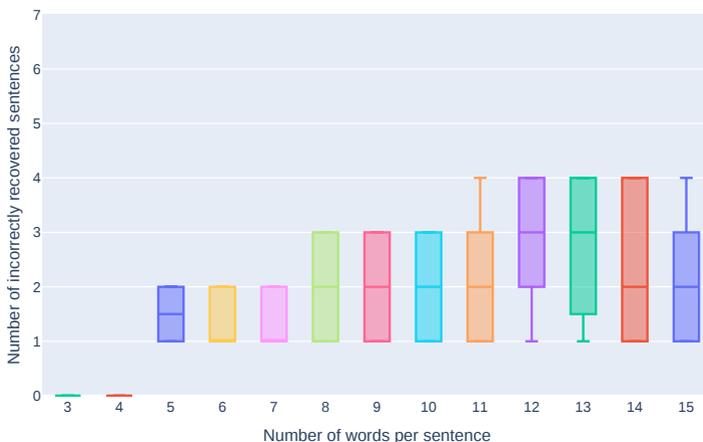


FIGURE 12. Numbers of incorrectly recovered sentences by Seq2Seq model from 10 test datasets

## 6. Conclusions, advantages and disadvantages of the proposed approach

The studies conducted in the work showed that the use of DNN models allows solving the problem posed in the work quite effectively. The results of text recovery from documents of the «Roscadastre» PLC acceptable for solving the practical problem are explained by the features of the dataset formation – it contains all pairs with incomplete and corresponding full texts. In addition, the incomplete text was considered in this work as a sequence of a certain number of adjacent words, which significantly simplified the process of its comparison with the full text. Sequences of arbitrary words of the full text were not included in the dataset.

The advantage of the proposed approach is the simplicity of the models

<b>Сведения об уточняемых земельных участках и их частях</b>
--

Сведения об уточняемых земельных участках и их частях

Сведения об уточняемых координатах, м

FIGURE 13. Correctly and incorrectly restored text "Information on clarified" (see Figure 1 and Figure 2)

TABLE 3. Metric values of Encoder-Decoder and Seq2Seq transformer models

Loss metric	Accuracy metric	BLEU metric	ROUGE-L metric
<b>Encoder-Decoder model</b>			
0.087	0.9798	0.3609	0.3071
<b>Seq2Seq transformer model</b>			
0.005	0.9982	0.5931	0.5173

and the features of the dataset formation for their training and validation. The disadvantage is the impossibility of restoring the text from a set of its arbitrary (inconsistent) words without significantly complicating the model, which involves analyzing the context of the sentence. The only reason for incorrect text recovery is related to the rather rare cases (1-3% of the total number) of the formation of identical embeddings (vectorization of incomplete text and the indicator of the document area in which the corresponding full text may be found, see section 3). An example of correct and incorrect recovery of incomplete text with identical embeddings is shown in Figure 13.

The results obtained in this work are planned for use in the information system of the «Roscadstr» PLC control and processing center for the purpose of converting scanned documents into their text analogues. To implement this process, the Seq2Seq transformer model was selected as it showed the best result compared to the Encoder-Decoder model, table. 3.

The metrics provided in this table are relevant only to the dataset created from the documents of the «Roscadastr» PLC. They may differ for datasets of other subject areas.

## References

- [1] N. C. Sabharwal, A. Agrawal. *Hands-on Question Answering Systems with BERT: Applications in Neural Networks and Natural Language Processing*, Apress, Berkeley, CA, 2021, ISBN 978-1-4842-6664-9, xv+184 pp. [doi](#) [↑76](#)
- [2] K. Aitken, V V. Ramasesh, Y. Cao, N. Maheswaranathan. *Understanding how encoder-decoder architectures attend*, 2021, 24 pp. [arXiv](#) [2110.15253](#) [cs.LG] [↑76](#)
- [3] A. Rahali, M. A. Akhloufi. "End-to-end transformer-based models in textual-based NLP", *Artificial Intelligence*, 4:1 (2023), pp. 54–110. [doi](#) [↑76](#)
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin. *Attention is all you need*, 2017, 15 pp. [arXiv](#) [1706.03762](#) [↑76](#), 87

- [5] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu. “BLEU: a method for automatic evaluation of machine translation”, *ACL’02 Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics* (July 7–12, 2002, Philadelphia, Pennsylvania, USA), ACL, Stroudsburg, 2002, pp. 311–318.  [↑76](#)
- [6] Ch.-Y. Lin. “ROUGE: a package for automatic evaluation of summaries”, *Proceedings of the Workshop on Text Summarization Branches Out*, WAS 2004 (July, 2004, Barcelona, Spain), ACL, 2004, 74–81 pp.  [↑76](#)
- [7] Vinokurov I. V.. “Using a convolutional neural network to recognize text elements in poor quality scanned images”, *Program Systems: Theory and Applications*, **13**:3(54) (2022), pp. 45–59.    [↑78](#)
- [8] Vinokurov I. V.. “Recognition of digital sequences using convolutional neural networks”, *Program Systems: Theory and Applications*, **14**:3 (2023), pp. 3–36 (in Russian, in English).    [↑78](#)
- [9] Th. Luong, H. Pham, Ch. D. Manning. “Effective approaches to attention-based neural machine translation”, *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (17–21 September, 2015, Lisbon, Portugal), ACL, 2015, ISBN 978-1-941643-32-7, pp. 1412–1421.   [↑79](#)
- [10] A. M. Dai, Q. V. Le. *Semi-supervised Sequence Learning*, NIPS 2015 (December 7–12, 2015, Montreal, Quebec, Canada), Advances in Neural Information Processing Systems, vol. **28**, Curran Associates, Inc., 2015, ISBN 9781510825024, 9 pp.   [↑79](#)
- [11] J. Gehring, M. Auli, D. Grangier, D. Yarats, Y. N. Dauphin. “Convolutional sequence to sequence learning”, *Proceedings of the 34th International Conference on Machine Learning* (6–11 August 2017, International Convention Centre, Sydney, Australia), PMLR, vol. **70**, 2017, pp. 1243–1252.   [↑79](#)
- [12] D. Ulyanov, A. Vedaldi, V. Lempitsky. “Deep image prior”, *International Journal of Computer Vision*, **128**:7 (2020), pp. 1867–1888.  [↑](#)
- [13] K. Hakala, A. Vesanto, N. Miekka, T. Salakoski, F. Ginter. *Leveraging text repetitions and denoising autoencoders in OCR post-correction*, 2019, 5 pp.   
 arXiv:1906.10907~[cs.CL] [↑79](#)
- [14] G. Huang, J. Wang, H. Tang, X. Ye. “BERT-based contextual semantic analysis for English preposition error correction”, *Journal of Physics: Conference Series*, **1693**:1 (2020), id. 012115, 5 pp.  [↑79](#)
- [15] K. Song, X. Tan, T. Qin, J. Lu, T.-Y. Liu. “MASS: masked sequence to sequence pre-training for language generation”, International Conference on Machine Learning (9–15 June 2019, Long Beach, California, USA), PMLR, vol. **97**, 2019, pp. 5926–5936.  arXiv:1905.02450~[cs.CL] [↑79](#)
- [16] Sh. Chollampatt, D. T. Hoang, H. T. Ng. “Adapting grammatical error correction based on the native language of writers with neural network joint models”, *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, EMNLP 2016 (1–4 November, 2016, Austin, Texas, USA), ACL, 2016, ISBN 978-1-945626-25-8, pp. 1901–1911.   [↑79](#)
- [17] A. Graves. *Supervised Sequence Labelling with Recurrent Neural Networks*, Studies in Computational Intelligence, vol. **385**, Springer, Berlin–Heidelberg, 2012, ISBN 978-3-642-24797-2, 146 pp.  [↑80](#)

- [18] T. Ge, X. Zhang, F. Wei, M. Zhou. “Automatic grammatical error correction for sequence-to-sequence text generation: an empirical study”, *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (July 28–August 2, 2019, Florence, Italy), ACL, 2019, ISBN 978-1-950737-48-2, pp. 6059–6064.    
- [19] X. Zhang, J. Zhao, Y. LeCun. “Character-level convolutional networks for text classification”, 2016, 9 pp. arXiv:1509.01626 [cs.LG]  
- [20] Z. Xie, A. Avati, N. Arivazhagan, D. Jurafsky, A. Ng. *Neural language correction with character-based attention*, 2016, 10 pp. arXiv:1603.09727 [cs.CL] 
- [21] J. Ramirez-Orta, E. Xamena, A. Maguitman, E. Milios, A. Soto. “Post-OCR document correction with large ensembles of character sequence-to-sequence models”, *Proceedings of the AAAI Conference on Artificial Intelligence*, **36** (2022), pp. 11192–11199.   
- [22] A. A. Alkhazraji, K. Baheerj, A. M. N. Alzubaidi. “Ancient textual restoration using deep neural networks: a literature review”, *2023 Al-Sadiq International Conference on Communication and Information Technology*, AICCIT 2023 (04–06 July 2023, Al-Muthana, Iraq), 2023, ISBN 9798350341898, pp. 64–69.  
- [23] F. Chollet. *Deep Learning with Python*, 2nd ed., Manning, 2021, ISBN 9781617296864, 504 pp.   
- [24] A. Géron. *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, 2nd ed., O’Reilly Media, Sebastopol, 2019, ISBN 978-1-492-03264-9, 848 pp.   
- [25] A. Kapoor, A. Gulli, S. Pal. *Deep Learning with TensorFlow and Keras: Build and deploy supervised, unsupervised, deep, and reinforcement learning models*, 3rd ed., Packt Publishing, 2022, ISBN 978-1803232911, 698 pp.  

Received	03.03.2024;
approved after reviewing	14.04.2024;
accepted for publication	15.08.2024;
published online	23.09.2024.

### Information about the author:



Igor Victorovich Vinokurov

Candidate of Technical Sciences (PhD), Associate Professor at the Financial University under the Government of the Russian Federation. Research interests: information systems, information technologies, data processing technologies

 0000-0001-8697-1032  
e-mail:

*The author declare no conflicts of interests.*

УДК 004.932.75'1, 004.89

10.25209/2079-3316-2024-15-3-75-110



## Восстановление текстовых последовательностей с использованием моделей глубокого обучения

Игорь Викторович **Винокуров**<sup>1✉</sup>

<sup>1</sup>Финансовый Университет при Правительстве Российской Федерации, Москва, Россия

<sup>✉</sup>[igvvinokurov@fa.ru](mailto:igvvinokurov@fa.ru)

**Аннотация.** В статье приведены результаты формирования, обучения и оценки качества работы моделей с архитектурами Encoder-Decoder и Sequence-To-Sequence (Seq2Seq) для решения задачи дополнения неполных текстов. Задачи такого типа достаточно часто возникают при восстановлении содержимого документов по их некачественным изображениям. Проведённые в работе исследования ориентированы на решение практической задачи формирования электронных копий отсканированных документов ППК «Роскадастр», распознавание которых стандартными средствами затруднительно или невозможно.

Формирование и исследование моделей осуществлялось на языке Python с использованием высокоуровневого API пакета Keras. С целью обучения и исследования моделей был сформирован набор данных, состоящий из нескольких тысяч пар. Каждая пара этого набора представляла собой неполный и соответствующий ему полный тексты. Для оценки качества работы моделей осуществлялось вычисление значений функции потерь loss и метрик accuracy, BLEU и ROUGE-L. Loss и accuracy позволили оценить эффективность моделей на уровне предсказания отдельных слов. Метрики BLEU и ROUGE-L использовались для оценки сходства между полными и восстановленными текстами. Полученные результаты показали, что обе модели Encoder-Decoder и Seq2Seq справляются с задачей восстановления текстовых последовательностей из их фиксированного множества, однако модель на основе трансформера Seq2Seq позволяет достичь лучших результатов по скорости и качеству обучения. (*Связанные тексты статьи на английском и на русском языках*)

**Ключевые слова и фразы:** модели глубокого обучения, Encoder-Decoder, трансформер Sequence-To-Sequence, восстановление текста, BLEU, ROUGE-L, Keras, Python

Для цитирования: Винокуров И. В. *Восстановление текстовых последовательностей с использованием моделей глубокого обучения* // Программные системы: теория и приложения. 2024. Т. 15. № 3(62). С. 75–110. (Англ.+русс.)  
[https://psta.psisras.ru/read/psta2024\\_3\\_75-110.pdf](https://psta.psisras.ru/read/psta2024_3_75-110.pdf)

## Введение

В последние годы с помощью моделей глубокого обучения (*Deep Neuaral Network*, DNN) получены существенные результаты в области обработке естественного текста (*Neural Language Processing*, NLP) [1]. Анализ литературных источников показал, что самыми распространёнными моделями, используемыми в таких задачах как преобразование (перевод) текста, восстановления текста из искаженных или непонятных документов, отсканированных документов плохого качества, нечитаемых рукописей, размытых или поврежденных изображений и т.п. являются Encoder-Decoder и трансформеры Seq2Seq.

Архитектура Encoder-Decoder, основанная на рекуррентных нейронных сетях (*Recurrent Neural Networks*, RNN) или свёрточных нейронных сетях (*Convolutional Neural Network*, CNN), состоит из двух основных компонентов – энкодера и декодера [2]. Энкодер преобразует входные данные во внутреннее представление, учитывая ключевые особенности их содержания. Декодер использует это представление для генерации выходных данных, последовательно предсказывая их элементы.

В отличие от этой модели, архитектура трансформера Seq2Seq [3] предлагает альтернативный подход к представлению последовательностей. Она использует механизм трансформации, основанный на множестве слоёв с механизмами внимания [4], что позволяет этой модели эффективно обрабатывать длинные текстовые последовательности. Механизм внимания позволяет модели учитывать важность отдельных слов в контексте всего предложения, способствуя тем самым генерации более качественно и связного текста. По сравнению с Encoder-Decoder, трансформер Seq2Seq обладает рядом преимуществ, таких как лучшая способность обращаться с длинными последовательностями, более гибкая архитектура и возможность обучения моделей на больших объёмах данных.

Для оценки качества моделей в NLP могут быть использованы как обычные метрики loss, ассугасу и т.п., так и метрики, специфичные для оценки качества сгенерированного текста, основными из которых являются BLEU (*Bilingual Evaluation Understudy*) [5] и ROUGE-L (*Recall-Oriented Understudy for Gisting Evaluation – Longest Common Subsequence*) [6]. Первая из них измеряет схожесть между предсказанным и эталонным текстом. Она использует синтаксическую информацию для сравнения последовательностей из  $n$  слов ( $n$ -грамм). Чем больше совпадений в  $n$ -граммах между предсказанным и эталонным текстами, тем выше будет значение

BLEU. Однако, эта метрика не учитывает семантическую и контекстную связь между словами, что может ограничивать её применимость. Вторая из метрик ROUGE-L оценивает качество автоматической суммаризации текста. Она сравнивает длину наибольшей общей последовательности слов между предсказанным и эталонным текстом с длиной эталонного текста, измеряя тем самым покрытие предсказанного текста относительно эталонного и позволяет оценить степень сжатия информации в сгенерированном тексте.

Основанием для проведения исследований, результаты которых приведены в этой статье, явилось невозможность восстановления текста на отсканированных документах плохого качества современными OCR-системами, рисунок 1.

Сведения об уточняемых земельных участках и их частях						
Сведения о характерных точках границы уточняемого земельного участка с кадастровым номером						
Обозначения характерных точек границы	Существующие координаты, м		Уточненные координаты, м		Средняя квадратическая погрешность положения характерных точек границ ЗУ, м	Описание закрепления точки
1	2	3	4	5	6	7

Рисунок 1. Документ с осветлёнными участками текста. Фрагменты текста, распознаваемые OCR, выделены цветом

Очевидным решением этой задачи является разработка простой системы соответствия предложений. Однако, как было отмечено выше, DNN-модели способны изучать сложные нелинейные зависимости между входными и выходными данными, что позволяет им более эффективно моделировать контекст и семантику текста. Кроме того, DNN-модели могут быть более гибкими и способными к обобщению, что делает их более эффективными при работе с различными типами текста и задачами восстановления информации. *Именно наличие у DNN-моделей обобщающих свойств послужило обоснованием для их использования при решении поставленной задачи – восстановленное и близкое по смыслу предложение лучше, чем его полное отсутствие.*

В разделе 1 осуществляется обоснование необходимости исследований и постановка задачи. Раздел 2 посвящён анализу работ по восстановлению

текстовых последовательностей. Создание набора данных для обучения моделей описано в разделе 3. Формирование и исследование моделей Encoder-Decoder и Seq2Seq приведено в разделе 4 и разделе 5 соответственно. Достоинства и недостатки восстановления текстовых последовательностей с использованием сформированных моделей приведены в разделе 6.

## 1. Постановка цели и задач исследования

При распознавании изображений текстовых документов стандартными средствами OCR не всегда можно получить их качественную копию в текстовом формате. Причиной являются размытые, неразборчивые или зашумлённые (например, в виде заметок от руки) участки текста на изображении документа [7, 8].

*Целью данной работы* является исследование применимости основных DNN-моделей для NLP – Encoder-Decoder и трансформера Seq2Seq для восстановления текста из заданного набора данных.

*Поставленная в работе цель может быть достигнута за счёт решения следующих основных задач:*

- (1) Формирование набора данных, заключающееся в подготовке пар в виде неполного и соответствующего ему полного текста.
- (2) Формирование оптимальных для достижения поставленной цели моделей Encoder-Decoder и трансформера Seq2Seq.
- (3) Анализ качества работы моделей в результате вычисления функции потерь и метрик ассигасы, BLEU и ROUGE-L.

## 2. Анализ работ по восстановлению текстовых последовательностей

Анализ доступных публикаций показал, что существует два основных типа моделей, которые могут быть использованы для сопоставления и восстановления текстовых последовательностей.

- (1) Модели с механизмом внимания (в большинстве случаев это модели трансформеров Seq2Seq), способные сфокусироваться на контексте предложения и выбрать наиболее подходящий вариант для замены или восстановления пропущенных слов.
- (2) Модели на основе RNN и CNN, позволяющие учитывать контекст предложения и, в случае CNN, особенности его визуализации на изображении.

Ниже приведено описание наиболее значимых, по мнению автора, работ по использованию моделей этих типов для восстановления оригинального текста.

В [9] авторы рассматривают и демонстрируют эффективность двух классов механизма внимания для сопоставления текстов на разных языках. Первый учитывает все исходные слова, второй рассматривает только подмножество исходных слов.

Восстановления пропущенных слов с использованием моделей на основе вариационных автоэнкодеров (*Variational Autoencoder*, VAE) предлагается, например, в [10]. Здесь VAE реализуют прогнозирование входной текстовой последовательности для последующего улучшения обучения RNN.

Модель Seq2Seq, позволяющая выполнить задачу генерации исправленного текста с использованием механизма внимания, описана в [11]. Модель может использоваться для преобразования неверной, некорректной последовательности символов в исправленный текст.

Модель на основе архитектуры Seq2Seq с механизмом внимания для исправления ошибок в тексте, полученном из OCR системы, приведена в [13]. Модель обучается на достаточно большом количестве текстовых пар, распознанных с помощью OCR.

Исследование применимости моделей BERT (*Bidirectional Encoder Representations from Transformers*) для исправления ошибок в предложениях неанглийских носителей языка, приведено в [14]. Модель способна использовать контекст предложения для исправления грамматических и стилистических ошибок в тексте.

Предварительно обученная модель с архитектурой маскированного трансформера Seq2Seq для восстановления фрагмента текста по его оставшейся части рассматривается в [15]. Кодер модели принимает на вход предложение со случайно замаскированным фрагментом (несколько последовательных токенов), а его декодер пытается предсказать этот замаскированный фрагмент. В работе показано, что в результате тонкой настройки модель способна достаточно точно восстанавливать исходных текст.

Авторы статьи [16] предлагают метод исправления ошибок в тексте с помощью нейронной сети, основанной на символьном самовнимании. Модель использует уровень символов для более точного исправления опечаток, орфографических и других типов ошибок в тексте.

В [17] автор описывает использование RNN для языковых моделей и предлагает методы для генерации текста в заданном контексте. Издание посвящено маркировке контролируемых последовательностей – важной области машинного обучения, включающей такие задачи как распознавание речи, рукописного ввода и маркировка частей речи.

Исследование применимости нейросетевых моделей для восстановления искажённых символьных изображений проводится в [18]. Авторы описывают методы восстановления испорченного текста, основанные на сочетании CNN и RNN.

В статье [19] представлена модель на основе CNN для автоматического исправления опечаток в тексте. Модель обучается на большом корпусе текстов и способна исправлять опечатки с высокой точностью.

В [20] приведён метод восстановления искажённых изображений с использованием CNN без необходимости обучения на большом наборе данных. Метод был изначально разработан для восстановления изображений, но может быть также применен и к восстановлению испорченного текста.

Авторы статьи [21] представляют архитектуру автоэнкодера для исправления ошибок в OCR-системах. Модель с автоэнкодером используется для извлечения скрытых представлений текста и их последующего восстановления.

Помимо указанных выше, достаточно большое количество работ по восстановлению текстовых последовательностей посвящено восстановлению текстов на разного типа изображениях исторических документов, повреждённых с течением времени. Их подробный анализ приведён в [22].

### 3. Формирование набора данных

Для обучения моделей Encoder-Decoder и Seq2Seq и последующего исследования их работы был сформирован собственный набор данных, представляющий собой, как уже отмечалось выше, множество пар вида:

$$\text{неполный текст} \rightarrow [\textit{start}]\text{полный текст}[\textit{end}]$$

Для формирования датасета использовались основные типы документов ППК «Роскадастр». Все уникальные предложения этих документов в сформированном датасете представляют собой полный текст. Результаты удаления из полного текста непрерывных последовательностей от 1

до  $N-3$  слов образует совокупность соответствующих ему неполных текстов (здесь  $N$  – количество слов в тексте). Удаление из текста именно непрерывных последовательных слов объясняется спецификой размытия и расположением неразборчивых участков текста на отсканированных документах, см. рисунок 1.

Количество предложений с полным текстом, входящих в сформированный датасет не превышает 250, количество соответствующих им неполных предложений, вошедших в датасет, составляет порядка 3000. Пример соответствия одного другому приведён на рисунке 2.

**Полный текст (электронный документ)**

Сведения об уточняемых земельных участках и их частях

**Неполный текст (результат распознавания OCR)**

1. Сведения об
2. об уточняемых
3. Сведения об уточняемых
4. Сведения об уточняемых земельных
5. об уточняемых земельных
6. Сведения об уточняемых земельных участках
7. уточняемых земельных участках
8. Сведения об уточняемых земельных участках и их
9. земельных участках и их частях
10. об уточняемых земельных участках и их
- ...

Рисунок 2. Полный текст и соответствующие ему первые 10 неполных текстов

Токенизация и векторизации текстовых пар осуществлялась с использованием `TextVectorization` из пакета `Keras`. Для снятия неоднозначности при восстановлении полных текстов, стандартизация текста перед токенизацией предполагала включение признака области документа, в которой он может присутствовать.

Для обучения DNN использовалось 80% пар из сформированного датасета, для валидации и тестирования – по 10%.

#### 4. Формирование и исследование модели Encoder-Decoder

Формирование и исследование модели Encoder-Decoder осуществлялось на языке Python с использованием API `Keras` [23–25]. На рисунке 3

приведена оптимальная структура модели, полученная в результате проведения экспериментальных исследований.

Энкодер этой модели состоит из следующих слоёв:

- (1) Входной слой (**InputLayer**), принимает входные данные в виде текстовых последовательностей.
- (2) Слой эмбединга (**Embedding**), преобразует слова текстовых последовательностей в их векторные представления.
- (3) Рекуррентный слой (**BidirectionalGRU**) обрабатывает последовательности и в прямом и обратном направлениях и выдаёт скрытое состояние (вектор контекста).

Слои декодера модели:

- (1) Входной слой (**InputLayer**) принимает входные данные в виде последовательности слов. Определяет форму и тип входных данных.
- (2) Слой эмбединга (**Embedding**) преобразует входные токены в плотные векторные представления заданной размерности, аналогично энкодеру.
- (3) Рекуррентный слой (**GRU**) принимает векторные представления и последовательно обрабатывает их, генерируя полный текст и учитывая контекст из энкодера. Инициализируется состоянием, сгенерированным энкодером.
- (4) Выходной слой (**Dense**) принимает выходные данные из декодера в текущем временном шаге и преобразует их в вектор вероятностей, каждая компонента которого относится к возможному следующему токenu восстановленного текста.

Описание основных параметров всех слоёв этой модели приведено в таблице 1.

Для оценки качества работы модели на уровне отдельных слов текстовых последовательностей осуществлялось вычисление значений функции потерь и метрики ассурасу. На рисунках 4 и 5 приведены значения функции потерь `sparse_categorical_crossentropy` и метрики ассурасу модели для 30-ти эпох обучения.

Количество эпох найдено экспериментальным путём и является оптимальным. При обучении модели использовался алгоритм стохастической оптимизации `rmsprop`.

Средние значения метрик BLEU и ROUGE-L, позволяющих оценить точность восстановления неполного текста из набора данных для тестирования, показаны на рисунке 6. Значения из диапазона 0.3-0.4 соответствуют пониманию и приемлемой трансляции текста. Для вычисления значений метрик использовались функции `sentence_bleu()` и `get_scores()` из пакетов

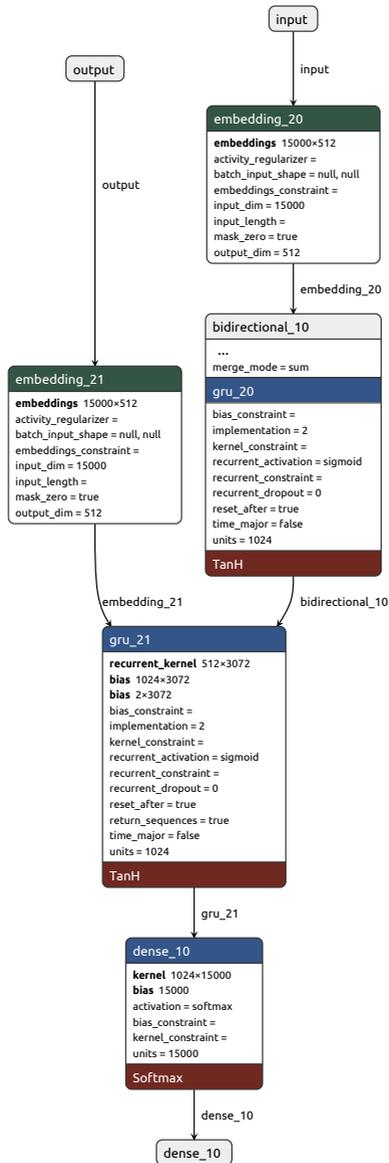


Рисунок 3. Модель Encoder-Decoder в API Keras для восстановления неполного текста

ТАБЛИЦА 1. Слои модели Encoder-Decoder

Тип слоя (имя на рисунке 3)	Функция активации	Входной тензор	Выходной тензор
<b>InputLayer</b> (input)	–	[(None, None)]	[(None, None)]
<b>Embedding</b> (embedding_20)	–	(None, None)	(None, None, 512)
<b>Bidirectional(GRU)</b> (bidirectional_10)	<code>tanh</code>	(None, None, 512)	(None, 1024)
<b>InputLayer</b> (output)	–	[(None, None)]	[(None, None)]
<b>Embedding</b> (embedding_21)	–	(None, None)	(None, None, 512)
<b>GRU</b> (gru_21)	<code>tanh</code>	[(None, None, 512), (None, 1024)]	(None, None, 1024)
<b>Dense</b> (dense_10)	<code>softmax</code>	(None, None, 1024)	(None, None, 15000)

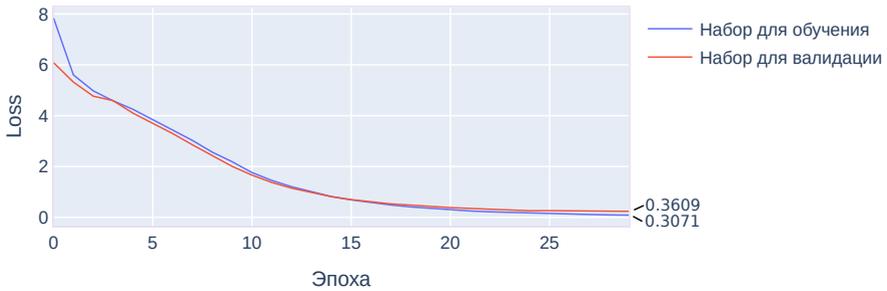


РИСУНОК 4. Потери модели

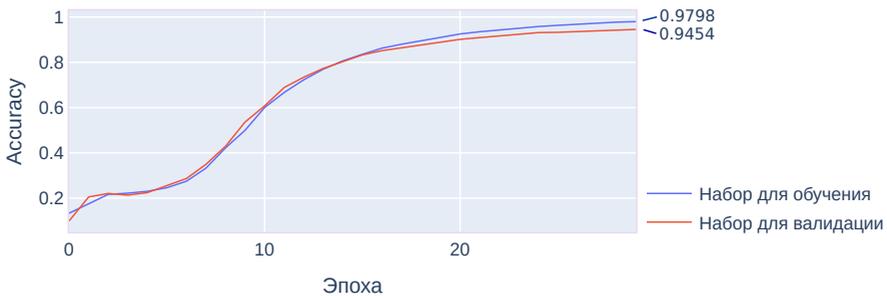


РИСУНОК 5. Точность модели

NLTK и Rouge соответственно. Их вызов осуществлялся после завершения очередной эпохи обучения из callback-функций метода `fit()`.



Рисунок 6. Метрики оценки качества восстановления текста

На рисунке 7 приведены результаты восстановления текста из 10 произвольно сформированных наборов данных по 250–300 предложений

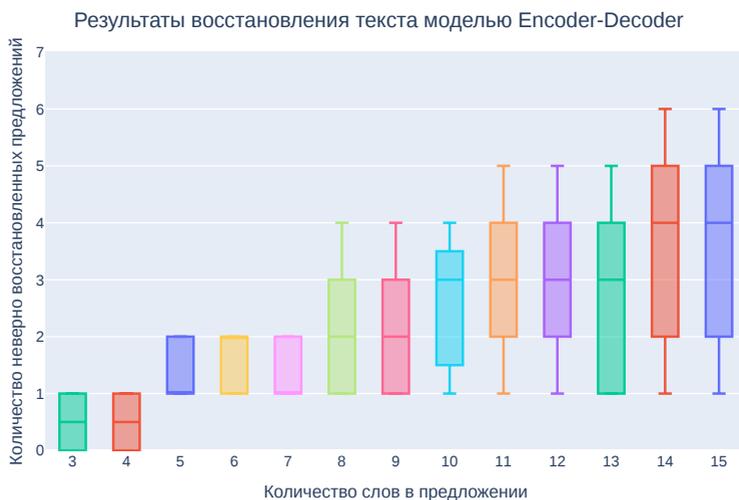


Рисунок 7. Количество неверно восстановленных предложений моделью Encoder-Decoder из 10-ти тестовых наборов данных

для тестирования. Количество неверно восстановленных предложений, состоящих из 11–15 слов, составило от 1 до 6.

## 5. Формирование и исследование модели трансформера Seq2Seq

Как и для предыдущей модели, формирование и исследование модели с архитектурой Seq2Seq осуществлялось на языке Python с использованием API Keras [23–25].

Оптимальная структура модели Seq2Seq, полученная в результате экспериментальных исследований, приведена на рисунке 8.

Модель трансформера состоит из следующих слоёв:

- (1) Входной слой (**InputLayer**) принимает входные данные в виде последовательности слов. Определяет форму и тип входных данных.
- (2) Слой позиционных эмбеддингов (**Positional Embedding Layer**) добавляет информацию о позиции слова в последовательности. Это позволяет модели учитывать порядок слов во входной и выходной последовательностях.
- (3) Слои энкодера трансформера (**TransformerEncoder Layer**), каждый из которых реализует механизм внимания и содержит полносвязные слои. За счёт этого они позволяют моделировать зависимости в последовательностях, извлекать признаки из входных данных и представлять их в оптимальном внутреннем представлении для более сложных вычислений и задач обработки естественного языка
- (4) Слои декодера трансформера (**TransformerDecoder Layer**). Аналогично энкодеру, декодер состоит из нескольких слоев декодера трансформера, которые также включают механизм внимания и полносвязные слои. Декодер генерирует выходную последовательность на основе контекстных представлений энкодера.
- (5) Выходной слой (**Dense**) преобразует предсказанные представления токенов в вероятностное распределение по всем возможным токенам.

Для создания более эффективных и связанных представлений текстовых последовательностей в слоях энкодера и декодера модели с архитектурой Seq2Seq реализован механизм внимания [4]. В таблице 2 приведено описание основных параметров всех слоёв этой модели.

ТАБЛИЦА 2. Слои модели трансформера Seq2Seq

Тип слоя (имя на рисунке 8)	Функция активации	Входной тензор	Выходной тензор
<b>InputLayer</b> (input)	–	[(None, None)]	[(None, None)]
<b>PositionalEmbedding</b> (positional_embedding)	–	(None, None)	(None, None, 256)
<b>TransformerEncoder</b> (transformer_encoder)	relu	(None, None, 256)	(None, None, 256)
<b>InputLayer</b> (output)	–	[(None, None)]	[(None, None)]
<b>PositionalEmbedding</b> (positional_embedding_1)	–	(None, None)	(None, None, 256)
<b>TransformerDecoder</b> (transformer_decoder)	relu	(None, None, 256)	(None, None, 256)
<b>Dense</b> (dense_4)	softmax	(None, None, 256)	(None, None, 15000)

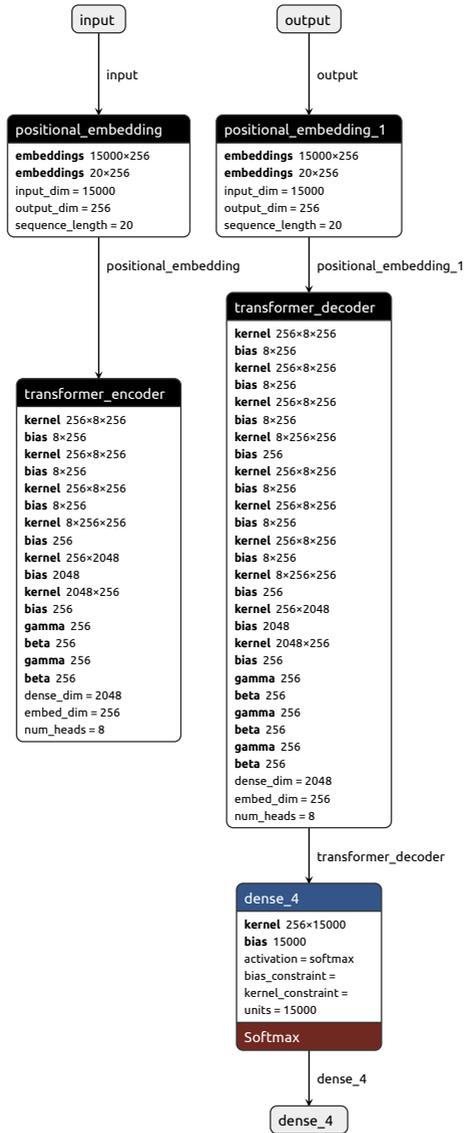


Рисунок 8. Эncoder и декодер модели трансформера Seq2Seq в API Keras для восстановления неполного текста

Исследование точности работы модели, по аналогии с предыдущей, заключалось в вычислении значений функции потерь и метрики точности – `sparse_categorical_crossentropy` и `accuracy` соответственно для каждой из 30-ти эпох её обучения, рисунок 9, 10. Количество эпох найдено экспериментальным путём и является оптимальным. При обучении модели использовался алгоритм стохастической оптимизации `gmsprop`.

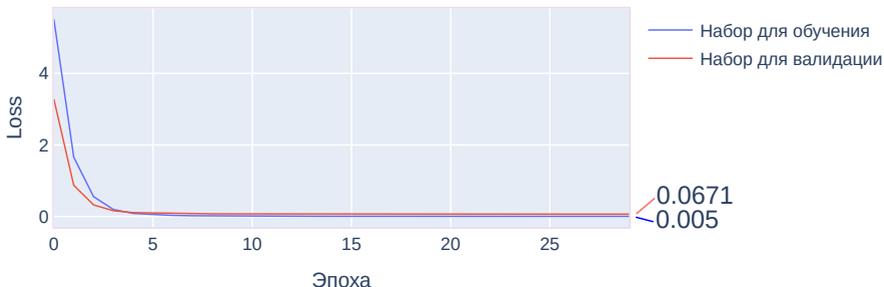


Рисунок 9. Потери модели

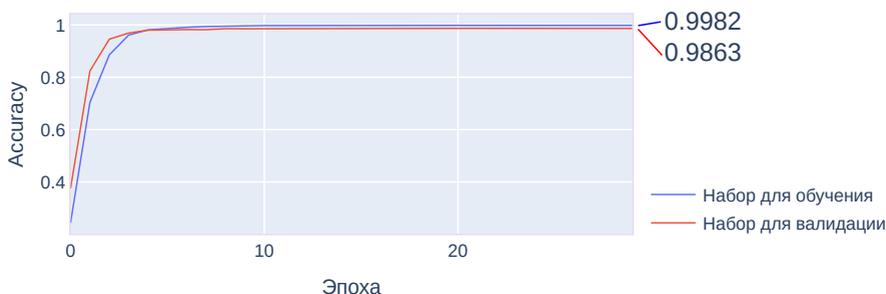


Рисунок 10. Точность модели

Средние значения метрик BLEU и ROUGE-L, позволяющих оценить точность восстановления моделью неполного текста из набора данных для тестирования, показаны на рисунке 11. Значения из диапазона 0.5-0.6 соответствуют высокому качеству трансляции текста.

На рисунке 12 приведены результаты восстановления неполного текста из 10-ти произвольно сформированных наборов данных для тестирования. Количество неверно восстановленных предложений меньше, чем для предыдущей модели, и составляет от 1 до 4. Как следствие, отметить вполне закономерный результат – модель трансформера Seq2Seq, за счёт использования механизма внимания, эффективнее модели Encoder-Decoder. Механизм внимания позволяет трансформеру выделять ключевые элементы входных и выходных последовательностей и более эффективно выявлять в них долгосрочные зависимости, что делает модель способной к обучению на сложных задачах генерации (в данном случае – восстановления) текста.

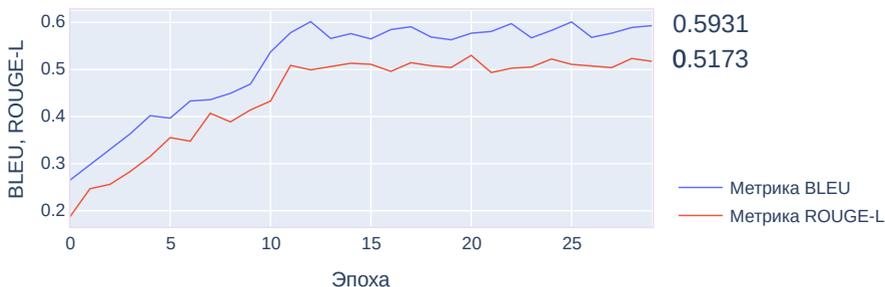


Рисунок 11. Метрики оценки качества восстановления текста

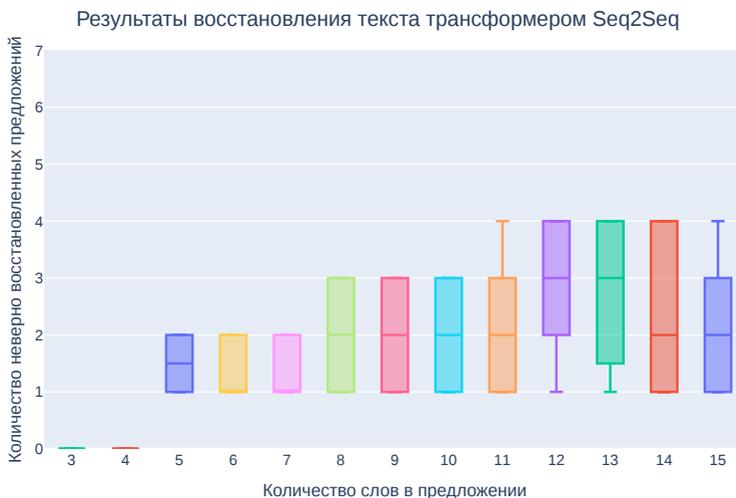


Рисунок 12. Количество неверно восстановленных предложений моделью Seq2Seq из 10-ти тестовых наборов данных

## 6. Выводы, достоинства и недостатки предложенного подхода

Проведённые в работе исследования показали, что использование DNN-моделей позволяет достаточно эффективно решить поставленную в работе задачу. Приемлемые для решения практической задачи результаты восстановления текста из документов ППК «Роскадастр» объясняются особенностями формирования набора данных – он содержит все пары с неполным и соответствующим ему полным текстами. Помимо этого, неполный текст рассматривался в данной работе как последовательность определённого количества соседних слов, что в значительной мере упрощало процесс его сопоставления с полным текстом. Последовательности из произвольных слов полного текста в набор данных не включались.

Достоинством предложенного подхода является простота моделей и особенностей формирования набора данных для их обучения и валидации. Недостаток – невозможность восстановления текста по совокупности его произвольных (непоследовательных) слов без существенного усложнения модели, предполагающей анализ контекста предложения. Единственная причина неправильного восстановления текста связана с достаточно редкими случаями (1-3% от общего количества) формирования одинаковых эмбеддингов (векторизация неполного текста и признак области документа, в котором возможно нахождение соответствующего ему полного текста, см. раздел 3). Пример правильного и неправильного восстановления неполного текста с одинаковыми эмбеддингами приведён на рисунке 13.

**Сведения об уточняемых земельных участках и их частях**

**Сведения об уточняемых земельных участках и их частях**

**Сведения об уточняемых координатах, м**

Рисунок 13. Правильно и неправильно восстановленный текст «Сведения об уточняемых» (см. рисунок 1 и 2)

Результаты, полученные в этой работе, планируются к использованию в информационной системе (ИС) ППК «Роскадстр» с целью преобразования отсканированных документов в их текстовые аналоги. Для реализации этого процесса выбрана модель трансформера Seq2Seq, как показавшая лучший по сравнению с моделью Encoder-Decoder результат, таблица 3.

Таблица 3. Значения метрик моделей Encoder-Decoder и Seq2Seq

Метрика loss	Метрика accuracy	Метрика BLEU	Метрика ROUGE-L
<b>Модель Encoder-Decoder</b>			
0.087	0.9798	0.3609	0.3071
<b>Модель трансформера Seq2Seq</b>			
0.005	0.9982	0.5931	0.5173

Метрики, приведённые в этой таблице, имеют отношение только к набору данных, сформированному из документов ППК «Роскадстр». Для наборов данных других предметных областей они могут отличаться.

### Список использованных источников

- [1] N. C. Sabharwal, A. Agrawal *Hands-on Question Answering Systems with BERT: Applications in Neural Networks and Natural Language Processing*.– Berkeley, CA: Apress.– 2021.– ISBN 978-1-4842-6664-9.– xv+184 pp.  [↑95](#)
- [2] K. Aitken, V V. Ramasesh, Y. Cao, N. Maheswaranathan *Understanding how encoder-decoder architectures attend*.– 2021.– 24 pp. arXiv  [2110.15253](#) [cs.LG] [↑95](#)

- [3] A. Rahali, M. A. Akhlooufi *End-to-end transformer-based models in textual-based NLP* // Artificial Intelligence.– 2023.– Vol. 4.– No. 1.– Pp. 54–110. doi ↑95
- [4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, I. Polosukhin *Attention is all you need.*– 2017.– 15 pp. arXiv:1706.03762 ↑95, 105
- [5] K. Papineni, S. Roukos, T. Ward, W.-J. Zhu *BLEU: a method for automatic evaluation of machine translation* // *ACL'02 Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics* (July 7–12, 2002, Philadelphia, Pennsylvania, USA), Stroudsburg: ACL.– 2002.– Pp. 311–318. doi ↑95
- [6] Ch.-Y. Lin *ROUGE: a package for automatic evaluation of summaries* // *Proceedings of the Workshop on Text Summarization Branches Out, WAS 2004* (July, 2004, Barcelona, Spain).– ACL.– 2004.– 74–81 pp. URL ↑95
- [7] И. В. Винокуров *Использование свёрточной нейронной сети для распознавания элементов текста на отсканированных изображениях плохого качества* // Программные системы: теория и приложения.– 2022.– Т. 13.– № 3(54).– С. 29–43. doi URL doi star ↑97
- [8] И. В. Винокуров *Распознавание цифровых последовательностей с использованием свёрточных нейронных сетей* // Программные системы: теория и приложения.– 2023.– Т. 14.– № 3(58).– С. 3–36 (русс.+англ.). doi URL doi star ↑97
- [9] Th. Luong, H. Pham, Ch. D. Manning *Effective approaches to attention-based neural machine translation* // *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing* (17–21 September, 2015, Lisbon, Portugal).– ACL.– 2015.– ISBN 978-1-941643-32-7.– Pp. 1412–1421. doi URL ↑98
- [10] A. M. Dai, Q. V. Le *Semi-supervised Sequence Learning*, NIPS 2015 (December 7–12, 2015, Montreal, Quebec, Canada), Advances in Neural Information Processing Systems.– Vol. 28.– Curran Associates, Inc.– 2015.– ISBN 9781510825024.– 9 pp. URL doi ↑98
- [11] J. Gehring, M. Auli, D. Grangier, D. Yarats, Y. N. Dauphin *Convolutional sequence to sequence learning* // *Proceedings of the 34th International Conference on Machine Learning* (6–11 August 2017, International Convention Centre, Sydney, Australia), PMLR.– vol. 70.– 2017.– Pp. 1243–1252. URL doi ↑98
- [12] D. Ulyanov, A. Vedaldi, V. Lempitsky *Deep image prior* // *International Journal of Computer Vision.*– 2020.– Vol. 128.– No. 7.– Pp. 1867–1888. doi ↑
- [13] K. Hakala, A. Vesanto, N. Miekka, T. Salakoski, F. Ginter *Leveraging text repetitions and denoising autoencoders in OCR post-correction.*– 2019.– 5 pp. arXiv:1906.10907~[cs.CL] ↑98
- [14] G. Huang, J. Wang, H. Tang, X. Ye *BERT-based contextual semantic analysis for English preposition error correction* // *Journal of Physics: Conference Series.*– 2020.– Vol. 1693.– No. 1.– id. 012115.– 5 pp. doi ↑98
- [15] K. Song, X. Tan, T. Qin, J. Lu, T.-Y. Liu *MASS: masked sequence to sequence pre-training for language generation*, International Conference on Machine Learning (9–15 June 2019, Long Beach, California, USA), PMLR.– vol. 97.– 2019.– Pp. 5926–5936. URL arXiv:1905.02450~[cs.CL] ↑98
- [16] Sh. Chollampatt, D. T. Hoang, H. T. Ng *Adapting grammatical error correction based on the native language of writers with neural network joint models* // *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016* (1–4 November, 2016, Austin, Texas, USA).– ACL.– 2016.– ISBN

- 978-1-945626-25-8.– Pp. 1901–1911.   ↑<sub>98</sub>
- [17] A. Graves *Supervised Sequence Labelling with Recurrent Neural Networks*, Studies in Computational Intelligence.– Vol. **385**.– Berlin–Heidelberg: Springer.– 2012.– ISBN 978-3-642-24797-2.– 146 pp.  ↑<sub>99</sub>
- [18] T. Ge, X. Zhang, F. Wei, M. Zhou *Automatic grammatical error correction for sequence-to-sequence text generation: an empirical study* // *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics* (July 28–August 2, 2019, Florence, Italy).– ACL.– 2019.– ISBN 978-1-950737-48-2.– Pp. 6059–6064.   ↑<sub>99</sub>
- [19] X. Zhang, J. Zhao, Y. LeCun *Character-level convolutional networks for text classification*.– 2016.– 9 pp. arXiv: 1509.01626~[cs.LG]  ↑<sub>99</sub>
- [20] Z. Xie, A. Avati, N. Arivazhagan, D. Jurafsky, A. Ng *Neural language correction with character-based attention*.– 2016.– 10 pp. arXiv: 1603.09727~[cs.CL] ↑<sub>99</sub>
- [21] J. Ramirez-Orta, E. Xamena, A. Maguitman, E. Milios, A. Soto *Post-OCR document correction with large ensembles of character sequence-to-sequence models* // *Proceedings of the AAAI Conference on Artificial Intelligence*.– 2022.– Vol. **36**.– Pp. 11192–11199.   ↑<sub>99</sub>
- [22] A. A. Alkhazraji, K. Baheeja, A. M. N. Alzubaidi *Ancient textual restoration using deep neural networks: a literature review* // *2023 Al-Sadiq International Conference on Communication and Information Technology, AICCIT 2023* (04–06 July 2023, Al-Muthana, Iraq).– 2023.– ISBN 9798350341898.– Pp. 64–69.  ↑<sub>99</sub>
- [23] F. Chollet *Deep Learning with Python*, 2nd ed..– Manning.– 2021.– ISBN 9781617296864.– 504 pp.  ↑<sub>100, 104</sub>
- [24] A. Géron *Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow*, 2nd ed..– Sebastopol: O’Reilly Media.– 2019.– ISBN 978-1-492-03264-9.– 848 pp.  ↑<sub>100, 104</sub>
- [25] A. Kapoor, A. Gulli, S. Pal *Deep Learning with TensorFlow and Keras: Build and deploy supervised, unsupervised, deep, and reinforcement learning models*, 3rd ed..– Packt Publishing.– 2022.– ISBN 978-1803232911.– 698 pp. ↑<sub>100, 104</sub>

Поступила в редакцию	03.03.2024;
одобрена после рецензирования	14.04.2024;
принята к публикации	15.08.2024;
опубликована онлайн	23.09.2024.

Рекомендовал к публикации

д.ф.-м.н. А. М. Елизаров

**Информация об авторе:****Игорь Викторович Винокуров**

Кандидат технических наук (PhD), ассоциированный профессор в Финансовом Университете при Правительстве Российской Федерации. Область научных интересов: информационные системы, информационные технологии, технологии обработки данных.



0000-0001-8697-1032

**e-mail:**

Декларация об отсутствии личной заинтересованности: *благополучие автора не зависит от результатов исследования.*



## АВТОРСКИЙ УКАЗАТЕЛЬ

**Винокуров Игорь Викторович**

*Восстановление текстовых последовательностей с использованием моделей глубокого обучения* ..... **75**, 112

**Новиков Николай Андреевич**

*Воспроизведение отклика графена на действие внешнего электрического поля с использованием модели сильно взаимодействующих ближайших соседей* **3**, 19

**Панферов Анатолий Дмитриевич**

*Воспроизведение отклика графена на действие внешнего электрического поля с использованием модели сильно взаимодействующих ближайших соседей* **3**, 19

**Смирнов Александр Владимирович**

*Применение Сиамских нейронных сетей для классификации биомассы растений по визуальному состоянию* ..... **53**, 72

**Степанов Дмитрий Николаевич**

*Математическое моделирование и исследование оптимальной конфигурации оптической стереосистемы, состоящей из двух плоских зеркал* ..... **23**, 48

**Тищенко Игорь Петрович**

*Математическое моделирование и исследование оптимальной конфигурации оптической стереосистемы, состоящей из двух плоских зеркал* ..... **23**, 48

**Ульянова Анастасия Алексеевна**

*Воспроизведение отклика графена на действие внешнего электрического поля с использованием модели сильно взаимодействующих ближайших соседей* **3**, 19

VOL. 15      ISSUE 3(62)      2024

## AUTHOR INDEX

**Novikov** Nikolay

*Simulation the response of graphene to an external electric field using the exact tight-binding model* ..... **3, 19**
**Panferov** Anatolii

*Simulation the response of graphene to an external electric field using the exact tight-binding model* ..... **3, 19**
**Smirnov** Alexander Vladimirovich

*Application of Siamese neural networks to classify plant biomass by visual state* **53, 72**
**Stepanov** Dmitry Nikolaevich

*Mathematical modeling and research of the optimal configuration of an optical stereo system consisting of two flat mirrors* ..... **23, 48**
**Tishchenko** Igor Petrovich

*Mathematical modeling and research of the optimal configuration of an optical stereo system consisting of two flat mirrors* ..... **23, 48**
**Ulyanova** Anastasiya

*Simulation the response of graphene to an external electric field using the exact tight-binding model* ..... **3, 19**
**Vinokurov** Igor Victorovich

*Recovering text sequences using deep learning models* ..... **75, 93**

VOL. 15 ISSUE 3(62) 2024

## CONTENTS

Research Article

COMPUTATIONAL SCIENCE

ANATOLII PANFEROV<sup>✉</sup>, NIKOLAY NOVIKOV, ANASTASIYA ULYANOVA. *Simulation the response of graphene to an external electric field using the exact tight-binding model (In Russ.)* ..... 3–19, **20–22**

Research Article

COMPUTATIONAL SCIENCE

DMITRY N. STEPANOV<sup>✉</sup>, IGOR P. TISHCHENKO. *Mathematical modeling and research of the optimal configuration of an optical stereo system consisting of two flat mirrors (In Russ.)* ..... 23–49, **50–53**

Research Article

ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

ALEXANDER V. SMIRNOV<sup>✉</sup>. *Application of Siamese neural networks to classify plant biomass by visual state (In Russ.)* ..... 53–72, **73–74**

Research Article

ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING

IGOR V. VINOKUROV<sup>✉</sup>. *Recovering text sequences using deep learning models (In Engl., In Russ.)* ..... **75–92**, 93–110

Author index ..... **114**

Чтобы сменить язык страницы, кликните, пожалуйста, флаг в верхнем углу

Авторский указатель ..... **113**

Содержание ..... **2**