

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ
РОССИЙСКОЙ ФЕДЕРАЦИИ

федеральное государственное автономное
образовательное учреждение высшего образования
«Самарский национальный исследовательский университет
имени академика С.П. Королева»
(Самарский университет)

Институт информатики и кибернетики
Кафедра технической кибернетики

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА

**Применение алгоритмов машинного обучения для задач
захвата движения человека на видеоизображении**

по программе бакалавриата по направлению подготовки
01.03.02 Прикладная математика и информатика,
профиль «Компьютерные науки»

Обучающийся _____ А.А. Сорока
(подпись)

Научный руководитель ВКР,
доцент, к.ф.-м.н.

_____ Д.А. Савельев
(подпись)

Нормоконтролёр _____ С.В. Суханов
(подпись)

Самара 2024

МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ
РОССИЙСКОЙ ФЕДЕРАЦИИ

федеральное государственное автономное
образовательное учреждение высшего образования
«Самарский национальный исследовательский университет
имени академика С.П. Королева»

Институт информатики и кибернетики
Кафедра технической кибернетики

УТВЕРЖДАЮ
Заведующий кафедрой

_____ А.В. Куприянов

«____» _____ 20____ г.

**ЗАДАНИЕ НА ВЫПУСКНУЮ КВАЛИФИКАЦИОННУЮ РАБОТУ
БАКАЛАВРА**

обучающемуся группы 6409-010302D *Сорока Александру Александровичу*

Тема ВКР: *Применение алгоритмов машинного обучения для задач захвата движения человека на видеоизображении*

утверждена приказом по университету от «08» 04 2021 г. № 333-ст.

Исходные данные: двумерное волновое уравнение, метод разделения переменных Фурье, признаки сходимости функциональных рядов, метод вычислительного эксперимента, метод конечных разностей.

Перечень вопросов, подлежащих разработке в ВКР:

- получение решения краевой задачи для волнового уравнения в виде бесконечного ряда Фурье при различных способах описания входного импульса;*
- получение оценки остатка ряда Фурье;*
- разработка программного средства для численного моделирования и исследования погрешности;*
- обеспечение контроля погрешности решения волнового уравнения;*

- разработка программы разностного решения волнового уравнения и использование программного средства численного моделирования с помощью рядов Фурье для тестирования программы разностного решения.

Руководитель ВКР
доцент

(подпись) Д.А. Савельев
« 15 » 02 2022 г.

Задание принял к исполнению

(подпись) А.А. Сорока
« 15 » 02 2022 г.

Примечание: Задание в данном случае печатается с двух сторон одного листа, пункты задания и все подписи будут находиться на одном листе.

РЕФЕРАТ

Выпускная квалификационная работа бакалавра: 56 с., 16 рисунков, 1 таблица, 25 источников.

Презентация: 12 слайдов Microsoft PowerPoint.

НЕЙРОННЫЕ СЕТИ, БЕЗМАРКЕРНЫЙ ЗАХВАТ ДВИЖЕНИЯ,
ОБРАБОТКА ИЗОБРАЖЕНИЙ, АНИМАЦИЯ, НЕЙРОСЕТЕВОЙ АНАЛИЗ
ДАННЫХ

Работа посвящена исследованию применимости алгоритмов машинного обучения для решения задач захвата движения человека на видеоизображении.

Часть работы посвящена обзору алгоритмов машинного обучения для задачи захвата движения.

В другой части рассматривается программная реализация алгоритма определения двумерных ключевых точек.

Также, в работе рассматриваются методы преобразования двумерных ключевых точек в трехмерные, а также анализируются полученные результаты.

СОДЕРЖАНИЕ

Введение.....	7
1 Постановка задачи и цель работы	9
1.1 Описание задачи захвата движения в общем смысле	9
1.1.1 Захват движений с помощью специальных маркеров	10
1.1.2 Безмаркерный захват движения	11
1.2 Задача захвата движения с точки зрения математики.....	11
1.3 Цель работы	12
2 Алгоритмы машинного обучения для захвата движения	13
2.1 Свёрточные нейронные сети.....	13
2.1.1 Свёрточные слои	14
2.1.2 Слои пулинга	14
2.1.3 Функции активации	16
2.1.4 Полносвязные слои	16
2.2 Рекуррентные нейронные сети	18
2.3 Сравнение CNN и RNN/LSTM для захвата движения	19
3 Программная реализация алгоритма определения двумерных ключевых точек.....	22
3.1 Подготовка датасета	22
3.2 Общий подход и архитектура модели.....	23
3.2.1 Backbone сети	23
3.2.2 RoI Pooling	26
3.2.3 Классификационные и регрессионные слои, блок ключевых точек	27
3.3 Оптимизатор Stochastic Gradient Descent (SGD).....	28
3.4 Расписание скорости обучения с использованием MultiStepLR.....	29
3.5 Функция потерь и процесс обучения	31
3.6 Мониторинг, валидация и регуляризация	32
4 Методы преобразования двумерных ключевых точек в трехмерные	35
4.1 Трёхмерная реконструкция по нескольким изображениям.....	35

4.2	Одиночное изображение с использованием глубины	36
4.3	Использование учебных данных с аннотацией глубины	39
4.4	Собственная реализация.....	39
4.4.1	Описание модели MiDaS.....	39
4.4.2	Преобразование 2D ключевых точек в 3D координаты.....	40
5	Анализ полученных результатов	42
5.1	Описание эксперимента	42
5.2	Результаты работы модели определения двумерных ключевых точек.	43
5.3	Результаты получения трёхмерных ключевых точек.....	48
	Заключение	52
	Список использованных источников	54

ВВЕДЕНИЕ

Современные исследования в области компьютерного зрения и машинного обучения актуальны благодаря широкому спектру применения этих технологий в различных индустриях, включая автоматизированное видеонаблюдение, интерактивные системы, спортивный анализ и реабилитацию после травм. Особенно значимыми становятся методы захвата и анализа движений человека, которые позволяют улучшить интерфейсы человеко-машинного взаимодействия и повысить точность биомеханических исследований.

Для анализа движений человека на видео часто используются алгоритмы машинного обучения, которые могут автоматически распознавать и классифицировать различные типы движений из видеоданных. Основным преимуществом этих алгоритмов является способность обучаться на больших объемах данных и адаптироваться к новым, ранее неизвестным условиям, что делает их идеально подходящими для задач компьютерного зрения.

В данной работе рассматривается задача разработки и апробации программного средства, основанного на методах машинного обучения, для захвата и анализа движения человека в видеопотоке.

Данная работа содержит пять разделов.

В первом разделе рассматривается постановка задачи, вводятся базовые определения, исследуется применение систем захвата движения в реальных задачах и ставится цель работы.

Во втором разделе рассматриваются алгоритмы машинного обучения для решения задачи захвата движения человека. Приводится обоснование выбора определенного алгоритма для реализации собственного решения.

В третьем разделе проводится описание программной реализации алгоритма определения двумерных ключевых точек. Подробно описывается модель сети, а также другие компоненты, используемые при обучении.

В четвертом разделе приведены методы преобразования двумерных ключевых точек в трехмерные.

И, наконец, в пятом разделе приведены результаты работы собственной модели определения ключевых точек. Также, провизуализирован результат совместной работы модели MiDaS с собственной моделью определения двумерных ключевых точек. Проведен анализ полученных результатов.

1 Постановка задачи и цель работы

1.1 Описание задачи захвата движения в общем смысле

Захват движения (Motion Capture, MoCap) — это технология, которая позволяет записывать движение объектов, в особенности человеческого тела, и применять эти данные для анимации моделей в 3D-пространстве.

При анимации видеоигр, фильмов и создании контента для виртуальной реальности, технологии захвата движения (рисунок 1) играют ключевую роль, добавляя реализм и интерактивность. В спортивных науках, анализ видеозаписей тренировок и соревнований позволяет тренерам и спортсменам улучшать технику и стратегии. В медицине, особенно в реабилитации, захват движений используется для оценки и коррекции походки пациентов, что критически важно для успешного восстановления.



Рисунок 1 – Технология Motion Capture

1.1.1 Захват движений с помощью специальных маркеров

Маркерный захват движения — это методика, при которой на теле актера или спортсмена закрепляются специальные маркеры, которые отслеживаются с помощью камер и других датчиков. Эти маркеры могут быть различных типов: отражающие, светящиеся или магнитные. В зависимости от системы захвата, маркеры отслеживаются камерами, которые расположены вокруг зоны захвата, и таким образом записывают точное положение маркеров в пространстве.

На первом этапе на тело испытуемого устанавливаются маркеры. Количество и расположение маркеров могут значительно варьироваться в зависимости от специфики задачи и системы захвата.

Прежде чем начать сессию захвата, необходимо калибровать систему. Это включает настройку камер или датчиков для оптимального захвата движения маркеров.

В процессе записи камеры считывают положение каждого маркера в пространстве. Данные с камер агрегируются и обрабатываются для создания трехмерной модели движения.

После записи сырые данные обрабатываются для построения анимационной модели. Это включает в себя удаление шумов, интерполяцию пропущенных данных и преобразование данных маркеров в координаты скелетной анимации.

Преимуществами захвата движения с помощью физических маркеров являются:

- 1) высокая точность и детализация захвата;
- 2) эффективность при сложных движениях, таких как акробатика или детальные мимические движения.

Но, такая система обладает и недостатками:

- 1) необходимость использования специального оборудования и пространства для захвата;

2) дискомфорт испытуемых из-за необходимости носить маркеры и ограничений, связанных с ними;

3) возможные ошибки из-за перекрытия маркеров или их потери в процессе движения.

1.1.2 Безмаркерный захват движения

Безмаркерный захват движения — это более современная и гибкая альтернатива, которая позволяет анализировать движения без прямого контакта с объектом захвата. Этот метод использует алгоритмы машинного зрения для анализа видеоизображений и выделения ключевых точек тела человека без физических маркеров.

В отличие от маркерных систем, безмаркерные системы требуют только видеозапись с одной или нескольких камер. С помощью алгоритмов компьютерного зрения происходит распознавание формы тела, определение положения суставов и ключевых точек. На основе полученных данных строится модель движения, которая может быть использована для анимации или анализа.

Преимущества:

1) удобство и доступность использования;

2) возможность анализа естественных движений в реальных условиях без искажений, вызванных наличием маркеров или специального снаряжения.

Недостатки:

1) меньшая точность и детализация по сравнению с маркерными системами;

2) зависимость от качества и условий съемки.

1.2 Задача захвата движения с точки зрения математики

Задача определения ключевых точек человеческого тела на видео может быть сформулирована как задача компьютерного зрения, в которой необходимо определить координаты предварительно определенных анатомических меток, таких как суставы, на изображении или в видеоряде. Математически это можно представить следующим образом.

Пусть I обозначает кадр из видеопотока, а $P = (p_1, p_2, \dots, p_n)$ – набор ключевых точек, которые необходимо определить. Каждая точка p_i описывается своими координатами на изображении (x_i, y_i) . Задача алгоритма машинного обучения – максимизировать [1] вероятность правильного определения этих координат, основываясь на обучающем наборе данных, содержащем аннотированные изображения.

Процесс обучения модели заключается в минимизации функции потерь, которая оценивает разницу между предсказанными алгоритмом координатами и истинными координатами точек на обучающих данных. Одним из популярных выборов для этой функции потерь является сумма квадратов разностей между предсказанными и истинными значениями координат:

$$L = \sum_{i=1}^n (x_i - \hat{x}_i)^2 + (y_i - \hat{y}_i)^2,$$

где \hat{x}_i и \hat{y}_i – предсказанные координаты точек, а x_i и y_i – истинные координаты.

1.3 Цель работы

Цель данной работы — исследовать применимость модели машинного обучения, способную с высокой точностью определять ключевые точки человеческого тела на изображениях и видео. Основные задачи включают улучшение точности захвата движения при различных условиях освещения, оптимизацию алгоритмов для работы в реальном времени на стандартном оборудовании и минимизацию нужды в ручной корректировке и аннотировании больших объемов данных.

Конечная цель исследования — не только создать эффективную техническую систему, но и продемонстрировать, как такие технологии могут быть интегрированы в реальные прикладные области, предоставляя значимую пользу в медицинских, спортивных и развлекательных приложениях.

2 Алгоритмы машинного обучения для захвата движения

2.1 Свёрточные нейронные сети

Convolutional Neural Networks (CNN) являются основным инструментом в современном компьютерном зрении и имеют значительное применение в анализе визуальных данных. В контексте захвата движения, CNN применяются для автоматического распознавания и отслеживания человеческих ключевых точек в последовательности видеокадров.

Принцип работы CNN заключается в автоматическом извлечении признаков из входных изображений посредством операций свёртки, что делает их идеальными для задач, требующих обработку больших и сложных визуальных данных без необходимости ручного задания характеристик. CNN состоят из нескольких типов слоёв (рисунок 2) [2]:

- 1) свёрточные слои (convolutional layers);
- 2) слои пулинга (pooling layers);
- 3) функции активации;
- 4) полносвязные слои (fully connected layers).

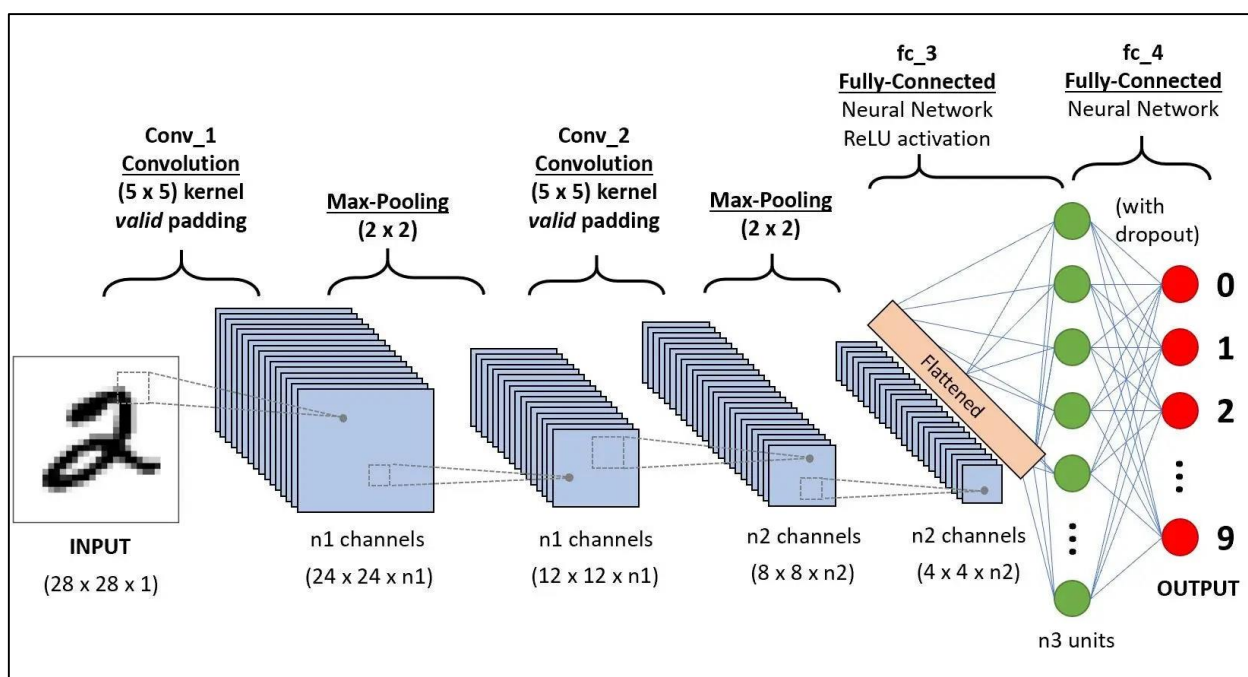


Рисунок 2 – Пример архитектуры свёрточной нейронной сети

2.1.1 Свёрточные слои

Свёрточные слои являются фундаментальной составляющей большинства архитектур глубокого обучения. Эти слои используют математическую операцию свертки для выделения важных признаков из входных данных.

Свёрточный слой использует набор фильтров или ядер, которые перемещаются по всему входному изображению, применяя математическую операцию свертки. Эта операция позволяет извлекать изображения из входных данных, создавая карты признаков, которые представляют собой агрегированную информацию о наличии определённых характеристик в различных регионах изображения [3].

В качестве входных данных выступают многоканальные матрицы данных, где каждый канал соответствует определенному цветовому каналу (например, RGB). Свертка применяется отдельно к каждому каналу, а результаты суммируются для получения итоговой карты признаков для каждого фильтра.

Математическая формулировка свёртки [4]: пусть I – входное изображение с размерами $H \times W \times D$, где D – количество каналов (например, для RGB $D = 3$). Фильтр свертки K имеет размеры $F \times F \times D$. Свертка вычисляется, как:

$$(I * K)(i, j) = \sum_{m=0}^{F-1} \sum_{n=0}^{F-1} \sum_{d=0}^{D-1} I(i + m, j + n, d) \cdot K(m, n, d),$$

где i, j – координаты элемента на выходной карте признаков, который получается после применения фильтра.

2.1.2 Слои пулинга

Слои пулинга, или слои субдискретизации, представляют собой ключевой компонент сверточных нейронных сетей, используемых для уменьшения размерности пространственных данных. Эти слои следуют непосредственно за сверточными слоями и играют важную роль в обеспечении инвариантности сети к масштабированию и другим искажениям изображения.

Суть пулинга заключается в применении операции уменьшения размерности к отдельным сегментам карт признаков, полученных после сверточных слоев. Это достигается путем применения агрегирующей функции, такой как максимум или среднее, к каждому такому сегменту.

Самый распространенный тип пулинга — максимальный пулинг [5]. В его рамках из каждого рассматриваемого подмножества входных данных выбирается максимальное значение. Например, если применять максимальный пулинг с размером окна 2×2 и шагом 2 к матрице признаков, каждое неперекрывающееся подокно 2×2 в этой матрице будет уменьшено до одного значения, равного максимальному из четырех.

Математически операция максимального пулинга для матрицы признаков A с размером окна $f \times f$ и шагом s можем быть описана следующим образом:

- 1) разделить матрицу A на неперекрывающиеся подматрицы размером $f \times f$;
- 2) для каждой подматрицы A_{sub} найти максимальное значение:

$$A'_{i,j} = \max(A_{sub}),$$

где i и j — индексы в результирующей матрице признаков A' , которая будет иметь уменьшенные пространственные размеры.

Представим, что у нас есть матрица признаков следующего вида:

$$A = \begin{bmatrix} 1 & 3 & 2 & 4 \\ 5 & 6 & 7 & 8 \\ 9 & 10 & 11 & 12 \\ 13 & 14 & 15 & 16 \end{bmatrix}$$

Применение операции максимального пулинга с размером окна 2×2 и шагом 2 приведет к следующей матрице:

$$A' = \begin{bmatrix} 6 & 8 \\ 14 & 16 \end{bmatrix}$$

Слои пулинга помогают уменьшить количество параметров и вычислений в сети, что способствует борьбе с переобучением и ускоряет обучение. Кроме того, благодаря уменьшению размерности, нейросеть становится менее

чувствительной к местоположению объектов во входном изображении, что улучшает её способность к обобщению [2].

2.1.3 Функции активации

Функции активации в нейронных сетях — это математические уравнения, которые определяют выходной сигнал нейрона на основе суммы входных сигналов. Эти функции необходимы для введения нелинейности в модель, что позволяет нейронной сети обучаться и выполнять более сложные задачи, чем просто линейная регрессия или классификация. Без нелинейности, какой бы сложной не была архитектура сети, она все равно оставалась бы линейной моделью, что значительно снижает её способность к обучению и аппроксимации функций.

Функция активации применяется к каждому нейрону в сети и определяет, насколько активным будет нейрон при данном входе. Проще говоря, это функция, которая включает или выключает нейрон. Рассмотрим этот процесс на примере одной из самых популярных функций активации — Rectified Linear Unit (ReLU) [6].

Формула для ReLU:

$$f(x) = \max(0, x)$$

Эта функция активации принимает один вход x и выдает x , если x положительное, и 0, если x отрицательное.

ReLU широко используется, потому что она проста в вычислении и помогает уменьшить вероятность исчезающего градиента, что часто встречается при использовании таких функций активации, как сигмоид или гиперболический тангенс. Важно отметить, что ReLU активирует нейроны только тогда, когда на входе есть активация, что делает нейронные сети разреженными, увеличивая тем самым эффективность и уменьшая вычислительные затраты.

2.1.4 Полносвязные слои

Полносвязные слои (dense layers) представляют собой основные строительные блоки в архитектуре нейронных сетей. В этих слоях каждый нейрон

предыдущего слоя соединён с каждым нейроном следующего слоя, что создаёт полную связность между слоями. Это позволяет сети интегрировать информацию, полученную на предыдущих этапах, для выполнения конкретных задач, таких как классификация или регрессия [7].

В контексте свёрточных нейронных сетей (CNN), полносвязные слои обычно располагаются в конце архитектуры после последовательности свёрточных и пулинг слоев. Основная функция этих слоев в CNN — синтезировать данные, извлечённые из предыдущих слоев, в предсказания, которые могут быть использованы для классификации или других типов вывода.

Каждый нейрон в полносвязном слое получает входы от всех активаций предыдущего слоя, умножает их на соответствующие веса, добавляет смещение (bias) и пропускает через функцию активации для получения выходного значения. Это можно математически представить следующим образом [8]:

$$y = f(Wx + b),$$

где x – вектор входных активаций из предыдущего слоя,

W – матрица весов,

b – вектор смещений,

f – функция активации,

y – вектор выходных активаций.

Рассмотрим полносвязный слой, который принимает входной вектор из трёх элементов и имеет два выходных нейрона. Предположим, что используется функция активации ReLU.

Инициализация:

$$x = [x_1, x_2, x_3],$$

$$W = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix},$$

$$b = [b_1, b_2]$$

Выход первого нейрона:

$$y_1 = \text{ReLU}(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b_1)$$

Выход второго нейрона:

$$y_2 = ReLU(w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + b_2)$$

Полносвязные слои широко используются во многих областях машинного обучения и глубокого обучения, включая:

- **Классификация изображений:** после извлечения признаков через свёрточные и пулинг слои, полносвязные слои используются для классификации изображений на основе этих признаков.
- **Регрессия:** в задачах, где требуется предсказать непрерывные значения, такие как цены на дома или температура, полносвязные слои могут обрабатывать признаки для предсказания этих значений.
- **Усиление признаков в сложных задачах:** в задачах с множеством переменных, таких как распознавание речи или машинный перевод, полносвязные слои помогают интегрировать и абстрагировать информацию на высоком уровне.

Эти слои обеспечивают возможность выражения сложных взаимосвязей между входными данными и желаемыми выходными значениями, что делает их неотъемлемой частью многих архитектур глубокого обучения.

2.2 Рекуррентные нейронные сети

Рекуррентные нейронные сети (рисунок 3) представляют собой класс нейросетевых моделей, специализированных на обработке последовательных данных, благодаря своей уникальной способности передавать информацию через временные шаги. Это достигается за счёт внутреннего состояния (или "памяти"), которое позволяет сети удерживать информацию о предыдущих данных в последовательности, делая RNN особенно подходящими для задач, где контекст имеет решающее значение, например, в анализе временных рядов, обработке естественного языка, синтезе речи и других.

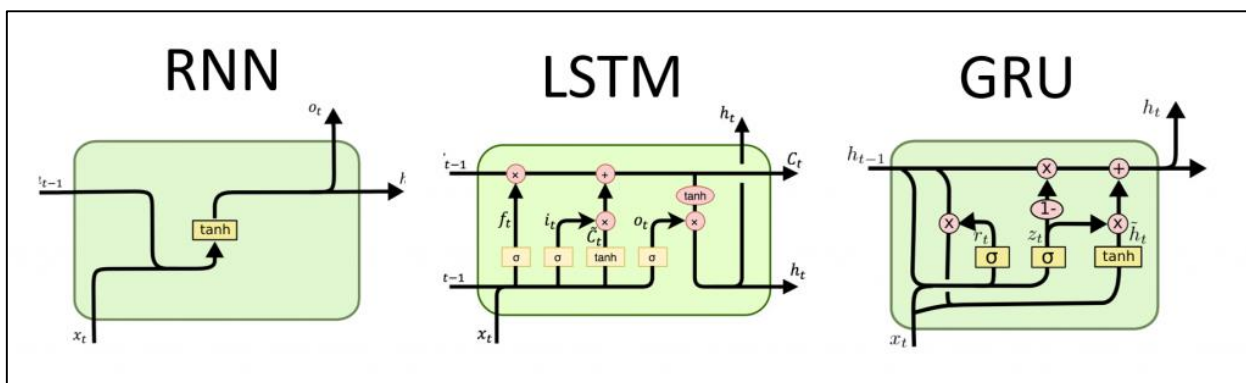


Рисунок 3 – Примерная архитектура рекуррентных нейронных сетей [9]

Одним из расширений RNN являются сети с долгой краткосрочной памятью (LSTM).

RNN и их вариации, включая LSTM и GRU, находят широкое применение во многих областях, где требуется анализ последовательностей данных:

- RNN могут анализировать исторические данные временных рядов, такие как акции или погодные данные, для предсказания будущих значений;
- в задачах NLP, таких как машинный перевод или распознавание речи, RNN могут улавливать семантические и синтаксические зависимости в тексте;
- используя последовательности нот и аккордов, RNN способны генерировать музыкальные произведения;
- RNN используются для понимания и предсказания действий и событий в видеопотоках, что может быть применено в системах видеонаблюдения или интерактивных мультимедийных системах.

2.3 Сравнение CNN и RNN/LSTM для захвата движения

Эффективность применения CNN и RNN/LSTM в этой задаче варьируется в зависимости от конкретных требований и условий применения.

Свёрточные нейронные сети идеально подходят для анализа визуальной информации благодаря своей способности эффективно обрабатывать пространственные отношения и текстуры на изображениях. В контексте MoCap, CNN могут быть использованы для точного определения ключевых точек и сегментов тела человека в различных позах и условиях освещения. Благодаря

своим свёрточным слоям, CNN способны автоматически выделять важные признаки из видеоданных, что значительно упрощает процесс обучения и уменьшает необходимость в ручном извлечении признаков и их предварительной обработке.

Однако, основным недостатком CNN является их неспособность учитывать временные зависимости в данных. В контексте захвата движения, где важно анализировать последовательность движений и каждое последующее состояние может зависеть от предыдущих, этот недостаток может существенно ограничивать применимость CNN.

В отличие от CNN, рекуррентные нейронные сети (RNN) и их усовершенствованный вариант — LSTM (Long Short-Term Memory) — разработаны специально для работы с последовательными данными. Эти сети способны обрабатывать информацию о предыдущих состояниях, что делает их идеальными для задач, где необходимо учитывать контекст, таких как анализ сложных движений и жестов в MoCap. LSTM улучшает возможности базовых RNN за счет механизмов забывания и сохранения информации, что позволяет эффективно управлять потоком информации и избегать проблемы исчезающих градиентов.

RNN и особенно LSTM требуют значительных вычислительных ресурсов для обучения, что может стать проблемой при работе с большими объемами данных, характерными для MoCap. Кроме того, настройка и оптимизация этих сетей может быть сложной из-за большого количества параметров и настроек, что увеличивает сложность моделей и время, необходимое для их тонкой настройки.

Выбор между CNN и RNN/LSTM для задач захвата движения должен базироваться на конкретных требованиях к задаче:

- 1) если важно пространственное распознавание и высокая точность в определении позы на отдельных кадрах, CNN будет предпочтительным выбором;

2) для анализа сложных движений, требующих учета временных зависимостей и последовательностей, более подходящим выбором окажется использование RNN или LSTM.

В идеальном сценарии комбинирование CNN для пространственного анализа и LSTM для обработки временной последовательности движений может дать наилучшие результаты. Такой подход позволит максимально использовать преимущества обеих технологий, обеспечивая точное распознавание поз и анализ сложных движений с учетом временной динамики.

3 Программная реализация алгоритма определения двумерных ключевых точек

3.1 Подготовка датасета

Датасеты для обучения и оценки моделей компьютерного зрения играют критически важную роль в разработке и улучшении алгоритмов, способных анализировать и интерпретировать визуальную информацию. Ключевым аспектом этих датасетов является их способность предоставлять разнообразные данные, которые отражают множество сценариев использования в реальном мире.

Современные датасеты включают изображения с различными уровнями сложности. Это может включать в себя разнообразные фоны, условия освещения и конфигурации объектов. Например, уличные сцены содержат множество перекрывающихся объектов, разные типы движения и разнообразные погодные условия. Датасеты типа ImageNet или COCO включают изображения, которые захватывают эти сложности, помогая моделям учиться на данных, максимально приближенных к реальности [10].

Для распознавания и локализации объектов важно, чтобы датасеты включали широкий спектр объектов различных форм, размеров и категорий. Например, датасет Pascal VOC включает объекты, такие как велосипеды, машины, люди и животные, каждый из которых имеет уникальные атрибуты и формы. Это важно для обучения моделей распознавать и точно классифицировать объекты в различных контекстах.

Хорошо аннотированные датасеты предоставляют точные и подробные метки, которые необходимы для обучения и оценки моделей. Это включает не только метки классов, но и ограничивающие рамки, аннотации ключевых точек и маски сегментации. Например, датасет COCO [11] предоставляет аннотации для сегментации на уровне пикселей и ключевые точки для анализа человеческой позы, что позволяет использовать его для задач, требующих детального понимания визуального контента.

Разнообразие в условиях съемки, таких как освещение, ракурс и качество изображения, также критически важно. Модели становятся устойчивыми к изменениям во входных данных, которые неизбежны в реальных приложениях. Датасеты, такие как ImageNet [12], включают изображения с различным разрешением и качеством, что помогает улучшить устойчивость и надежность обученных моделей.

Датасеты должны отражать реальные сценарии применения, чтобы обеспечить практическую значимость разработанных моделей. Например, датасеты для автономного вождения, такие как KITTI, предоставляют данные, собранные с датчиков на автомобилях, что помогает в разработке систем, способных работать в реальных дорожных условиях.

В рамках данной работы был использован датасет COCO (Common Objects in Context), который является одним из стандартов в индустрии для задач детекции, сегментации и распознавания ключевых точек. COCO содержит более 330 тысяч изображений с более чем 200 тысячами меток, включающих ограничивающие рамки, сегментационные маски и ключевые точки.

Кроме того, используется аугментации данных и дополнительные техники предварительной обработки для улучшения обучения и обобщающей способности модели. Это включает в себя масштабирование, повороты и изменение цветовых настроек изображений, что позволяет модели лучше адаптироваться к разнообразным условиям и улучшить ее способность к обобщению на новые данные.

3.2 Общий подход и архитектура модели

3.2.1 Backbone сети

Backbone сеть в архитектуре Keypoint R-CNN является фундаментальным компонентом, который отвечает за первичное извлечение признаков из входных изображений. В контексте глубокого обучения, backbone — это обычно предварительно обученная сверточная нейронная сеть (CNN), которая преобразует исходные изображения в сложный набор признаков или карт признаков. Эти карты признаков служат основой для всех последующих этапов

анализа и обработки, таких как детекция объектов, классификация и определение ключевых точек.

Распространенные архитектуры, используемые в качестве backbone в задачах компьютерного зрения, включают [13]:

– ResNet (Residual Networks): Эта сеть использует так называемые "остаточные блоки", которые помогают обучать очень глубокие нейронные сети. Основное математическое выражение для слоя в ResNet выглядит следующим образом:

$$x_{l+1} = x_l + F(x_l, W_l),$$

где x_l – входной вектор признаков на уровне l ,

F – остаточная функция,

W_l – веса, которые необходимо обучить [13].

– VGG (Visual Geometry Group): Простая и мощная архитектура, основанная на повторении блоков, состоящих из сверточных слоев с маленьким размером фильтра (3x3), за которыми следуют слои пулинга. В VGG все слои используют одинаковый шаг и дополнение нулями, что позволяет сохранять пространственные размеры через слои [14].

Feature Pyramid Network (FPN) является дополнением к стандартному backbone и предназначена для улучшения детекции объектов различного размера. FPN создает иерархию признаков на разных уровнях разрешения, что позволяет модели лучше адаптироваться к объектам разных масштабов. В FPN каждый уровень пирамиды создается путем слияния информации из двух источников: нижнего уровня с высоким разрешением и верхнего уровня с более глубокой семантической информацией. Это достигается с помощью операций свертки и апсемплинга:

$$P_l = U_l(P_{l+1}) \times C_l,$$

где P_l – признаки на уровне l пирамиды,

C_l – выходные данные сверточного слоя на уровне l ,

U_l – операция апсемплинга [15].

В архитектуре Keypoint R-CNN используется ResNet50 в качестве backbone сети, основываясь на её эффективности (таблица 1) и универсальности в различных задачах компьютерного зрения. ResNet50, с её остаточными блоками, способствует упрощению обучения глубоких сетей, предотвращая исчезновение градиента, что критически важно для обеспечения стабильности и высокой производительности модели.

Таблица 1 - Эффективность оценки ключевых точек на валидационном наборе COCO [13]

Backbone	AP	AP_{50}	AP_{75}	AP_M	AP_L
ResNet-50	70,4	88,6	78,3	67,1	77,2
Res2Net-50	71,5	89,0	79,3	68,2	78,4
ResNet-101	71,4	89,3	79,3	68,1	78,1
Res2Net-101	72,2	89,4	79,8	68,9	79,2
Res2Net-vgg-50	72,2	89,5	79,7	58,5	79,4
Res2Net-vgg-101	73,0	89,5	80,3	69,5	80,0

Таблица 1 представляет сравнение различных архитектур backbone на валидационном наборе данных COCO для оценки эффективности определения ключевых точек. Оценки включают общую точность (AP), точность при IoU пороге 0.5 (AP_{50}) и 0.75 (AP_{75}), а также точность для средних (AP_M) и больших (AP_L) объектов.

ResNet-50 показывает средние результаты с общей точностью 70,4, что ниже, чем у более сложных моделей. Однако, данная модель уже интегрирована в архитектуру Keypoint R-CNN как *keypointrcnn_resnet50_fpn* в среде PyTorch, что делает её простой в использовании для начальных этапов разработки и экспериментов, несмотря на то, что она не обеспечивает максимально возможной точности.

Res2Net-50 и Res2Net-101 обеспечивают лучшую общую точность, чем их аналоги ResNet, со значениями 71,5 и 72,2 соответственно. Это указывает на то, что более сложные модификации ResNet архитектуры могут лучше справляться с задачами определения ключевых точек благодаря улучшенному извлечению признаков и более глубокой обработке информации.

Res2Net-v1b-50 и Res2Net-v1b-101 показывают наивысшие результаты во всех категориях, с *AP* достигающим 73,0 для Res2Net-v1b-101. Это подчеркивает преимущества использования очень глубоких сетей с расширенной архитектурой для сложных задач определения ключевых точек.

3.2.2 RoI Pooling

RoI (Region of Interest) Pooling — это техника, используемая в сверточных нейронных сетях для извлечения фиксированного размера признаков из произвольно размерных регионов. Эта операция критически важна в задачах, где необходимо обработать локальные области изображения, такие как детекция объектов и сегментация. RoI Pooling был популяризирован архитектурами, такими как Fast R-CNN, для улучшения производительности моделей в этих задачах [16].

На рисунке 4 можно посмотреть на результат операции определения RoI.

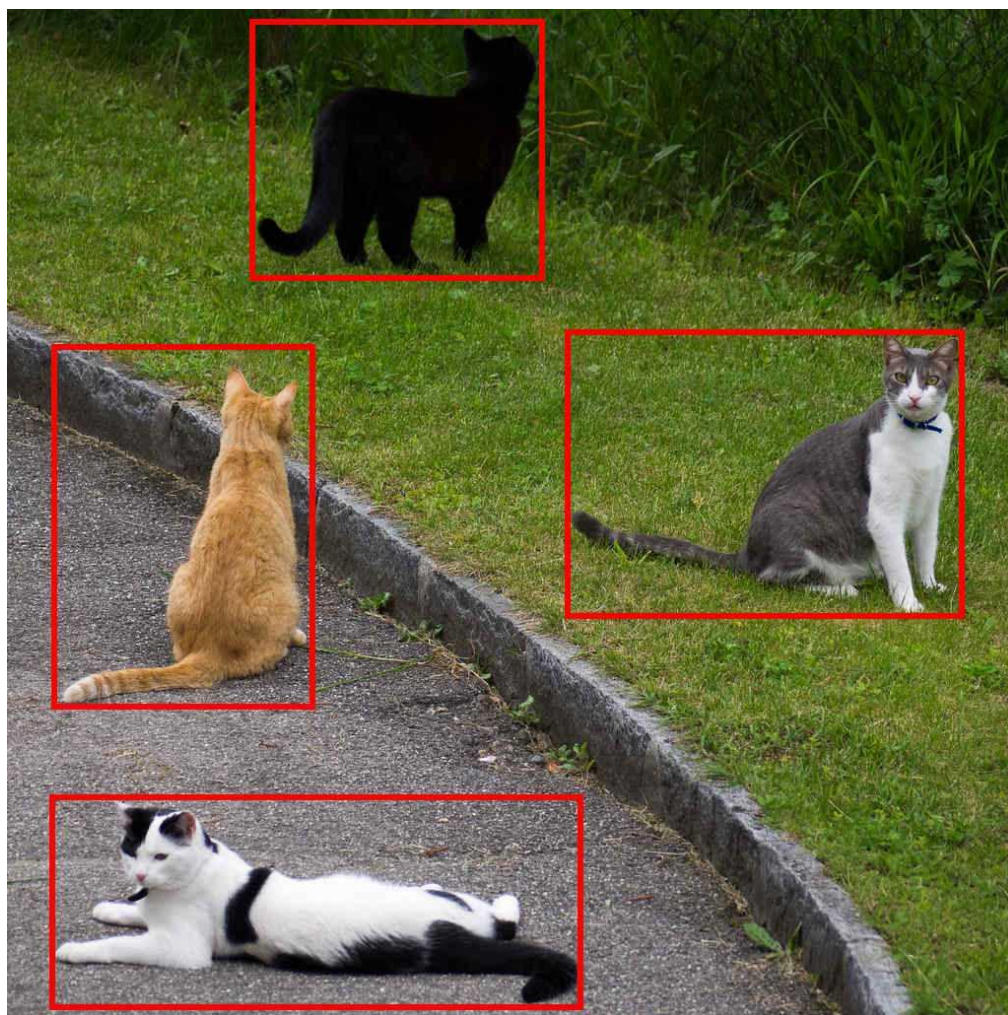


Рисунок 4 – Визуализация результатов работы определения RoI

Рассмотрим входную карту признаков F размером $H \times W \times D$, где H , W и D обозначают высоту, ширину и глубину карты признаков соответственно. Пусть дано N регионов интереса, каждый из которых задан четырьмя координатами: (x_1, y_1, x_2, y_2) , где (x_1, y_1) и (x_2, y_2) обозначают верхний левый и нижний правый углы региона на карте признаков.

Для каждого региона R_n RoI Pooling операция выполняется следующим образом [17]:

- 1) регион R_n преобразуется к фиксированному размеру $H' \times W'$, который задается заранее. Это делается для стандартизации выходных данных, позволяя использовать их на следующих слоях нейронной сети;
- 2) преобразованный регион делится на $H' \times W'$ ячеек равного размера;
- 3) в каждой ячейке выполняется операция пулинга для получения одного значения признака. Пулинг определяется, как [17]:

$$P(i, j) = \max_{(x, y) \in S_{ij}} F(x, y),$$

где $P(i, j)$ – это значение признака в ячейке (i, j) ,

S_{ij} – множество пикселей в ячейке (i, j) входной карты признаков.

3.2.3 Классификационные и регрессионные слои, блок ключевых точек

После извлечения признаков с помощью RoI Pooling, следующий шаг заключается в использовании классификационных и регрессионных слоев для определения класса объектов и точного позиционирования их ограничивающих рамок. Классификационный слой применяет функцию softmax для оценки вероятности принадлежности каждого предложенного региона к одному из возможных классов объектов:

$$p(c|r) = \frac{e^{s_c}}{\sum_{c'=1}^C e^{s_{c'}}},$$

где $p(c|r)$ – вероятность класса c для региона r ,

s_c – сырые оценки классификации для класса c ,

а C – общее количество классов.

Регрессионный слой затем адаптирует ограничивающие рамки, предложенные RPN, для того чтобы они как можно точнее обрамляли детектируемый объект. Это достигается путем корректировки четырех параметров рамки: центра x, y , ширины w и высоты h . Преобразование координат рамки выражается через:

$$(x, y, w, h) \leftarrow (x + \delta x, y + \delta y, we^{\delta w}, he^{\delta h}),$$

где $(\delta x, \delta y, \delta w, \delta h)$ - предсказания регрессии относительно исходной рамки [18].

После классификации объектов и уточнения их рамок, выполняется задача определения ключевых точек. Каждая детектированная рамка объекта снабжается набором признаков, из которых затем предсказываются ключевые точки. Это достигается с помощью регрессии для каждой точки внутри рамки:

$$k_{x,y} = (x + \Delta x, y + \Delta y),$$

где $k_{x,y}$ – координаты предсказанной ключевой точки,

(x, y) – исходные координаты точки в рамке,

$(\Delta x, \Delta y)$ – смещения, предсказанные сетью.

3.3 Оптимизатор Stochastic Gradient Descent (SGD)

Оптимизатор Stochastic Gradient Descent (SGD) — это один из основных методов оптимизации, используемых в обучении нейронных сетей и других вычислительных задачах, связанных с минимизацией функций. SGD модифицирует традиционный метод градиентного спуска путём обновления параметров модели на основе оценки градиента, полученной из случайно выбранного подмножества всего набора данных.

Формула для обновления весов в SGD:

$$w_{t+1} = w_t - \eta \nabla L(w_t),$$

где w_t – параметры модели на шаге t ,

η – скорость обучения, положительный коэффициент, определяющий размер шага в пространстве параметров,

$\nabla L(w_t)$ – градиент функции потерь L по параметрам w на шаге t , полученный из текущего батча данных.

В классическом градиентном спуске (Gradient Descent, GD) для обновления параметров модели используется весь набор данных для вычисления градиента функции потерь. Это требует значительных вычислительных ресурсов и времени, особенно при работе с большими данными. В отличие от GD, SGD обновляет параметры, используя только один обучающий пример за раз. Это делает SGD значительно быстрее, особенно в контекстах, где данные имеют высокую избыточность [19].

SGD идеально подходит для задач, где данные поступают последовательно, например, в системах рекомендаций или для персонализированных приложений, где модель должна адаптироваться к новым пользователям или меняющимся предпочтениям пользователя в реальном времени. Это обеспечивается возможностью модели обновляться постоянно, по мере поступления каждого нового данных, без необходимости повторного обучения с нуля.

Помимо достоинств, SGD обладает и рядом недостатков. Проблемы сходимости — один из основных недостатков SGD. Из-за того, что градиенты вычисляются исходя из одного образца, результаты могут сильно колебаться между итерациями. Это может привести к тому, что процесс обучения будет менее стабильным и потребует большего количества итераций для достижения сходимости, особенно если выбранный размер шага обучения не оптимален.

Риск застревания в локальных минимумах или седловых точках особенно актуален в контекстах, где функция потерь не является выпуклой. В таких случаях SGD может «застрять», не найдя глобальный минимум ошибки, который оптимален для всех данных. Это связано с тем, что локальные шумы и несовершенства данных могут ввести модель в заблуждение, указывая на неправильное направление обновлений.

3.4 Расписание скорости обучения с использованием MultiStepLR

MultiStepLR — это метод адаптации скорости обучения в процессе тренировки нейронной сети, реализуемый в библиотеке *PyTorch* через модуль

torch.optim.lr_scheduler. Этот метод позволяет планомерно уменьшать скорость обучения в predetermined моменты тренировки, что часто приводит к более эффективному и стабильному обучению моделей глубокого обучения.

Скорость обучения — один из наиболее важных гиперпараметров при тренировке нейронных сетей. Она определяет величину шага, с которым обновляются веса модели в направлении антиградиента функции потерь. Слишком большая скорость обучения может привести к тому, что обучение будет "перепрыгивать" через минимумы функции потерь, в то время как слишком маленькая скорость замедлит процесс обучения и может привести к застреванию в локальных минимумах.

MultiStepLR работает, уменьшая скорость обучения на заданный коэффициент γ в заданные моменты времени, что обычно совпадает с эпохами. Эти «шаги» или «ступеньки» уменьшения скорости обучения задаются списком эпох. Когда текущая эпоха обучения достигает одной из указанных эпох в списке, скорость обучения умножается на коэффициент γ .

Преимущества использования *MultiStepLR* включают возможность настройки контрольных точек для изменения скорости обучения, что позволяет адаптировать процесс обучения под конкретные этапы развития модели. Это особенно важно, поскольку разные фазы тренировки могут требовать различной скорости обучения для оптимальной эффективности. Кроме того, использование *MultiStepLR* способствует преодолению плато в обучении, когда прогресс модели замедляется и стандартная постоянная скорость обучения не приносит желаемых результатов. Постепенное уменьшение скорости обучения может также улучшить сходимость, позволяя модели более стабильно и глубоко аппроксимировать оптимальные значения весов, что снижает риск осцилляций вокруг минимума функции потерь.

Однако использование *MultiStepLR* не лишено недостатков. Одной из основных проблем является необходимость ручной настройки моментов изменения скорости и коэффициента уменьшения, что может потребовать значи-

тельных усилий и времени на эксперименты, особенно в начале обучения. Эффективность конкретного расписания скорости обучения может также значительно варьироваться в зависимости от специфики задачи и используемых данных, что ограничивает универсальность данного подхода. Кроме того, существует риск застревания в локальных минимумах или на седловых точках, если скорость обучения будет снижена слишком рано или если начальная скорость обучения была установлена слишком низкой, что может препятствовать достижению глобального минимума функции потерь [20].

3.5 Функция потерь и процесс обучения

Функция потерь, используемая в данной модели, включает в себя несколько компонентов, предназначенных для различных аспектов задачи. Основные используемые функции потерь:

– Smooth L1 Loss - это вариант функции потерь, который обычно используется для задач регрессии, таких как определение ограничивающих рамок в задачах детекции объектов. Эта функция потерь менее чувствительна к выбросам по сравнению с традиционной L2 потерей благодаря своей способности уменьшать влияние больших ошибок на процесс обучения. Формула Smooth L1 Loss представлена как [21]:

$$L_s(x) = \begin{cases} 0,5x^2, & \text{if } |x| < 1 \\ |x| - 0,5, & \text{otherwise} \end{cases},$$

где x – разница между предсказанным значением и истинным значением.

– Cross-Entropy Loss - часто применяется в задачах классификации. Эта функция измеряет различие между двумя вероятностными распределениями: предсказанным и истинным. Для двухклассовой классификации функция потерь вычисляется как:

$$L_c(p, y) = -(y \log(p) + (1 - y) \log(1 - p)),$$

где p – предсказанная вероятность принадлежности объекта к классу, y – истинная метка класса (0 или 1).

Процесс обучения модели разделяется на несколько ключевых этапов:

- 1) **прямое распространение:** на этом этапе входные данные подаются в модель, проходят последовательно через все слои сети, в результате чего модель генерирует предсказание. Каждый слой сети преобразует входные данные согласно своим обученным параметрам и передает результат следующему слою;
- 2) **вычисление потерь:** после получения предсказаний на выходе модели рассчитывается значение функции потерь. Это значение представляет собой оценку того, насколько предсказания модели отличаются от истинных значений;
- 3) **обратное распространение:** с помощью этого метода производится расчет градиентов функции потерь по всем параметрам модели — от выходных слоев к входным. Этот процесс помогает определить вклад каждого параметра в ошибку и определить, как необходимо изменить параметры, чтобы минимизировать функцию потерь;
- 4) **оптимизация:** на основе вычисленных градиентов и с использованием метода SGD параметры модели обновляются таким образом, чтобы функция потерь уменьшалась. Оптимизация является итеративным процессом, и каждое обновление параметров приближает модель к оптимальному решению задачи.

3.6 Мониторинг, валидация и регуляризация

Одним из основных этапов в процессе обучения является валидация модели. Валидация (рисунок 5) позволяет оценить, как модель будет работать в реальных условиях, используя тестовый набор данных, который не участвовал в обучении. Этот процесс включает в себя периодическую оценку модели на этом наборе данных, чтобы мониторить такие метрики, как точность классификации, точность локализации объектов (в случае задач детекции или сегментации), а также другие метрики, специфичные для конкретной задачи, такие как площадь под ROC-кривой (AUC) или среднее значение точности пересечения по объединению (IoU). Значения IoU после первой эпохи можно увидеть на рисунке 6.

Test: [2200/5000]	eta: 0:04:01	model_time: 0.0635 (0.0615)	evaluator_time: 0.0050 (0.0057)	time: 0.0731	data: 0.0010	max mem: 7163
Test: [2300/5000]	eta: 0:03:50	model_time: 0.0475 (0.0615)	evaluator_time: 0.0040 (0.0057)	time: 0.0606	data: 0.0009	max mem: 7163
Test: [2400/5000]	eta: 0:03:39	model_time: 0.0515 (0.0612)	evaluator_time: 0.0040 (0.0056)	time: 0.0635	data: 0.0009	max mem: 7163
Test: [2500/5000]	eta: 0:03:29	model_time: 0.0490 (0.0610)	evaluator_time: 0.0040 (0.0057)	time: 0.0596	data: 0.0009	max mem: 7163
Test: [2600/5000]	eta: 0:03:18	model_time: 0.0500 (0.0608)	evaluator_time: 0.0030 (0.0056)	time: 0.0597	data: 0.0011	max mem: 7163
Test: [2700/5000]	eta: 0:03:08	model_time: 0.0480 (0.0605)	evaluator_time: 0.0040 (0.0056)	time: 0.0578	data: 0.0007	max mem: 7163
Test: [2800/5000]	eta: 0:02:59	model_time: 0.0655 (0.0607)	evaluator_time: 0.0040 (0.0056)	time: 0.0907	data: 0.0010	max mem: 7163
Test: [2900/5000]	eta: 0:02:50	model_time: 0.0540 (0.0606)	evaluator_time: 0.0040 (0.0056)	time: 0.0638	data: 0.0009	max mem: 7163
Test: [3000/5000]	eta: 0:02:41	model_time: 0.0545 (0.0606)	evaluator_time: 0.0040 (0.0056)	time: 0.0674	data: 0.0009	max mem: 7163
Test: [3100/5000]	eta: 0:02:32	model_time: 0.0480 (0.0605)	evaluator_time: 0.0040 (0.0056)	time: 0.0634	data: 0.0010	max mem: 7163
Test: [3200/5000]	eta: 0:02:23	model_time: 0.0582 (0.0606)	evaluator_time: 0.0040 (0.0056)	time: 0.0661	data: 0.0009	max mem: 7163
Test: [3300/5000]	eta: 0:02:15	model_time: 0.0520 (0.0608)	evaluator_time: 0.0040 (0.0056)	time: 0.0838	data: 0.0009	max mem: 7163
Test: [3400/5000]	eta: 0:02:07	model_time: 0.0455 (0.0607)	evaluator_time: 0.0040 (0.0056)	time: 0.0552	data: 0.0012	max mem: 7163
Test: [3500/5000]	eta: 0:01:58	model_time: 0.0440 (0.0606)	evaluator_time: 0.0040 (0.0056)	time: 0.0555	data: 0.0010	max mem: 7163
Test: [3600/5000]	eta: 0:01:49	model_time: 0.0520 (0.0604)	evaluator_time: 0.0040 (0.0056)	time: 0.0699	data: 0.0009	max mem: 7163
Test: [3700/5000]	eta: 0:01:41	model_time: 0.0505 (0.0604)	evaluator_time: 0.0040 (0.0056)	time: 0.0668	data: 0.0011	max mem: 7163
Test: [3800/5000]	eta: 0:01:33	model_time: 0.0540 (0.0604)	evaluator_time: 0.0050 (0.0056)	time: 0.0795	data: 0.0011	max mem: 7163
Test: [3900/5000]	eta: 0:01:25	model_time: 0.0510 (0.0603)	evaluator_time: 0.0040 (0.0056)	time: 0.0572	data: 0.0010	max mem: 7163
Test: [4000/5000]	eta: 0:01:17	model_time: 0.0520 (0.0602)	evaluator_time: 0.0040 (0.0055)	time: 0.0659	data: 0.0010	max mem: 7163
Test: [4100/5000]	eta: 0:01:09	model_time: 0.0565 (0.0601)	evaluator_time: 0.0040 (0.0055)	time: 0.0626	data: 0.0009	max mem: 7163
Test: [4200/5000]	eta: 0:01:01	model_time: 0.0630 (0.0603)	evaluator_time: 0.0040 (0.0055)	time: 0.0881	data: 0.0011	max mem: 7163
Test: [4300/5000]	eta: 0:00:53	model_time: 0.0540 (0.0604)	evaluator_time: 0.0040 (0.0055)	time: 0.0596	data: 0.0010	max mem: 7163
Test: [4400/5000]	eta: 0:00:46	model_time: 0.0535 (0.0606)	evaluator_time: 0.0040 (0.0055)	time: 0.0703	data: 0.0010	max mem: 7163
Test: [4500/5000]	eta: 0:00:38	model_time: 0.0630 (0.0608)	evaluator_time: 0.0040 (0.0055)	time: 0.0983	data: 0.0011	max mem: 7163
Test: [4600/5000]	eta: 0:00:30	model_time: 0.0535 (0.0608)	evaluator_time: 0.0040 (0.0055)	time: 0.0638	data: 0.0012	max mem: 7163
Test: [4700/5000]	eta: 0:00:22	model_time: 0.0535 (0.0608)	evaluator_time: 0.0040 (0.0055)	time: 0.0675	data: 0.0010	max mem: 7163
Test: [4800/5000]	eta: 0:00:15	model_time: 0.0520 (0.0607)	evaluator_time: 0.0040 (0.0055)	time: 0.0670	data: 0.0007	max mem: 7163
Test: [4900/5000]	eta: 0:00:07	model_time: 0.0555 (0.0607)	evaluator_time: 0.0050 (0.0055)	time: 0.0704	data: 0.0011	max mem: 7163
Test: [4999/5000]	eta: 0:00:00	model_time: 0.0491 (0.0607)	evaluator_time: 0.0040 (0.0055)	time: 0.0600	data: 0.0007	max mem: 7163

Рисунок 5 – Валидация модели

IoU metric: keypoints						
Average Precision	(AP)	@[IoU=0.50:0.95	area=	all	maxDets= 20]	= 0.491
Average Precision	(AP)	@[IoU=0.50	area=	all	maxDets= 20]	= 0.753
Average Precision	(AP)	@[IoU=0.75	area=	all	maxDets= 20]	= 0.528
Average Precision	(AP)	@[IoU=0.50:0.95	area=medium	maxDets= 20]		= 0.462
Average Precision	(AP)	@[IoU=0.50:0.95	area= large	maxDets= 20]		= 0.559
Average Recall	(AR)	@[IoU=0.50:0.95	area=	all	maxDets= 20]	= 0.591
Average Recall	(AR)	@[IoU=0.50	area=	all	maxDets= 20]	= 0.842
Average Recall	(AR)	@[IoU=0.75	area=	all	maxDets= 20]	= 0.633
Average Recall	(AR)	@[IoU=0.50:0.95	area=medium	maxDets= 20]		= 0.543
Average Recall	(AR)	@[IoU=0.50:0.95	area= large	maxDets= 20]		= 0.660

Рисунок 6 – IoU метрики

Проведение регулярных проверок в процессе обучения помогает выявить такие проблемы, как переобучение, когда модель хорошо работает на тренировочных данных, но плохо на тестовых. Это также позволяет разработчикам принять меры для корректировки параметров обучения, выбора модели или самого процесса обучения в реальном времени.

Для предотвращения переобучения часто используются методы регуляризации. Одним из самых популярных методов является L2 регуляризация

[22], известная также как Weight Decay. Этот метод добавляет к функции потерь член, пропорциональный квадрату нормы вектора весов модели:

$$L = L_0 + \lambda \|w\|^2,$$

где L_0 – исходная функция потерь,

w – вектор весов модели,

λ – коэффициент регуляризации, который контролирует степень влияния штрафа на итоговую функцию потерь.

Этот штраф за большие веса помогает предотвратить переобучение, делая модель менее чувствительной к небольшим колебаниям во входных данных, тем самым способствуя лучшей обобщающей способности.

Также широко используется регуляризация Dropout, при которой случайным образом исключается часть нейронов в процессе обучения, что помогает уменьшить зависимость модели от конкретных атрибутов входных данных и тем самым увеличить её способность к обобщению.

4 Методы преобразования двумерных ключевых точек в трехмерные

4.1 Трёхмерная реконструкция по нескольким изображениям

Мультивидовая стереоскопия (рисунок 7) является одним из наиболее эффективных подходов для создания детализированных трехмерных моделей из двухмерных изображений. Процесс включает в себя анализ нескольких снимков объекта, сделанных с разных ракурсов, что позволяет воссоздать объемные структуры с высокой точностью [23].

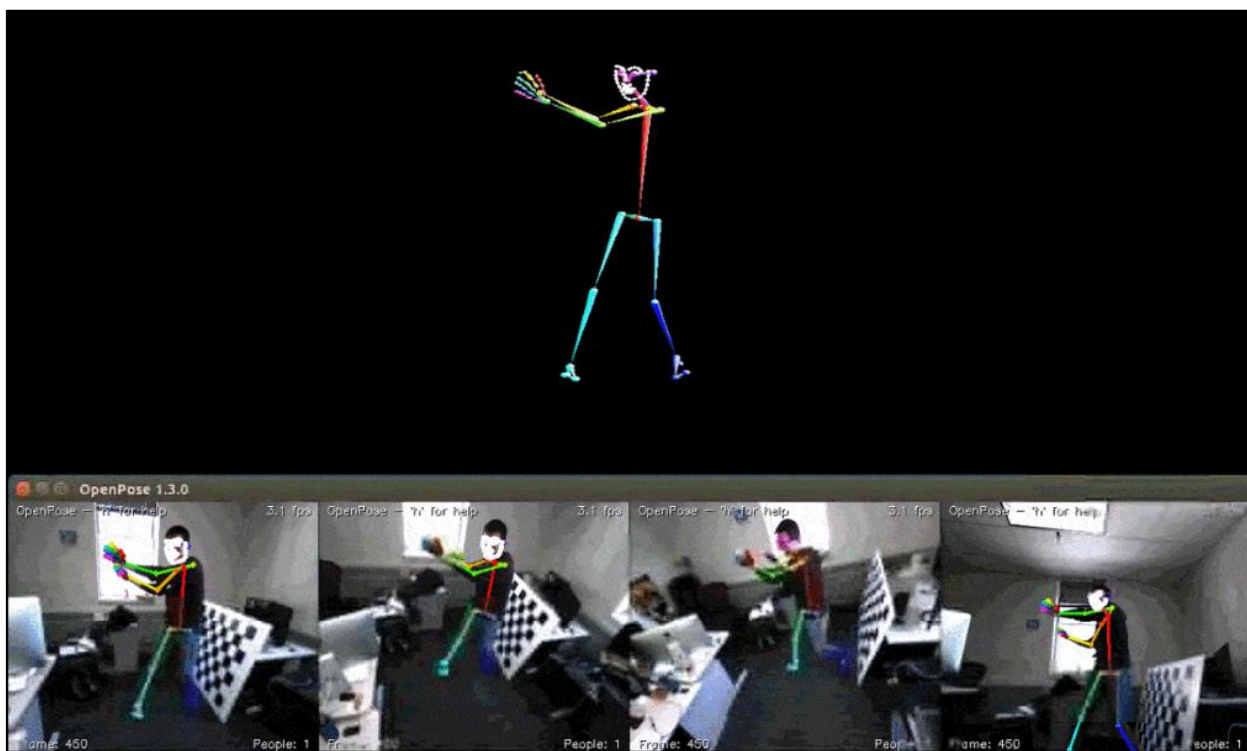


Рисунок 7 – Реконструкция 3D скелета из нескольких изображений

На начальном этапе происходит сбор данных. Данный этап является критически важным этапом, поскольку качество и количество изображений напрямую влияют на точность и детализацию трехмерной модели. Изображения должны быть сняты так, чтобы между ними было достаточное перекрытие, что позволит алгоритмам сопоставления точек находить соответствия с высокой точностью. Обычно используются специализированные камеры, расположенные под разными углами относительно объекта.

Затем, алгоритмы компьютерного зрения, такие как SIFT (Scale-Invariant Feature Transform) или SURF (Speeded Up Robust Features), используются для обнаружения уникальных особенностей на изображениях. Эти особенности должны быть инвариантны к масштабированию, вращению и изменениям освещения, чтобы обеспечить их надежное сопоставление между разными снимками.

Финальным этапом является триангуляция. Триангуляция - это процесс определения положения точки в трехмерном пространстве по её проекциям на несколько изображений. Для каждой сопоставленной точки строится система уравнений, основанная на параметрах проекции каждой камеры и координатах точек на изображениях. Решение этой системы позволяет найти координаты точки в глобальной системе координат. Математически этот процесс можно описать через минимизацию ошибки воспроизведения:

$$\min_X \sum_{i=1}^n \|x_i - P_i X\|^2,$$

где x_i – проекция точки X на изображение i ,

P_i – матрица проекции камеры i .

Мультивидовая стереоскопия широко используется в киноиндустрии для создания реалистичных CGI (computer-generated imagery) сцен, в археологии для реконструкции артефактов, в робототехнике для создания точных карт окружающей среды, и в медицине для создания трехмерных моделей анатомических структур.

4.2 Одиночное изображение с использованием глубины

Использование одиночного изображения с картой глубины актуально, когда доступ к мультимедийным настройкам ограничен или невозможен. Этот метод позволяет эффективно восстанавливать пространственную структуру объектов и сцен из одного изображения, используя информацию о глубине каждой точки [24].

Первым этапом является получение карты глубины. Основная идея данного подхода состоит в том, что глубина пикселя на изображении известна. Но вот возможностей ее получения несколько:

– Стереокамеры: используют две камеры (или больше) для захвата двух видов одной и той же сцены под немного разными углами. Глубина рассчитывается на основе различий между этими двумя видами, опираясь на принципы триангуляции;

– LiDAR: системы, основанные на LiDAR (рисунок 8), излучают лазерные лучи и измеряют время, необходимое им для возвращения после отражения от объектов. Это время преобразуется в расстояние, создавая высокоточные карты глубины;

– современные методы глубокого обучения могут предсказывать карту глубины (рисунок 9) непосредственно из одиночного изображения на основе обучения модели на больших наборах данных с аннотированной глубиной.

Следующим этапом является преобразование 2D ключевых точек в 3D. Используя карту глубины, преобразование 2D пиксельных координат изображения в 3D пространственные координаты осуществляется через процесс, известный как проекция обратного вида (back-projection). Этот процесс учитывает внутреннюю калибровку камеры, включая фокусное расстояние и центр проекции, для перевода пиксельных координат в реальные координаты. Математически это описывается следующими уравнениями:

$$x = (u - c_x) \cdot \frac{d}{f},$$

$$y = (v - c_y) \cdot \frac{d}{f},$$

$$z = d,$$

где u, v – координаты пикселя на изображении,

c_x, c_y – координаты центра проекции,

f – фокусное расстояние камеры,

d – значение глубины для пикселя.

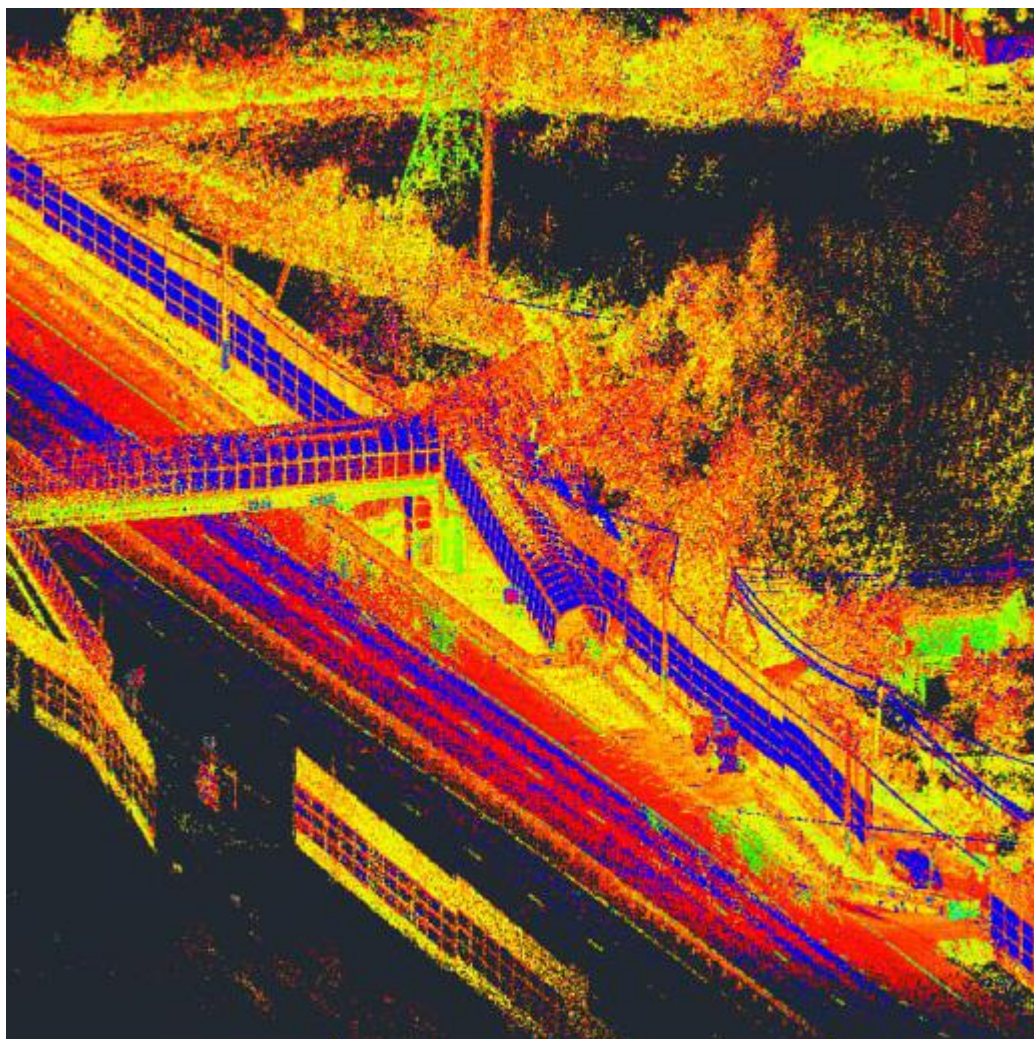


Рисунок 8 – Лидарная съемка

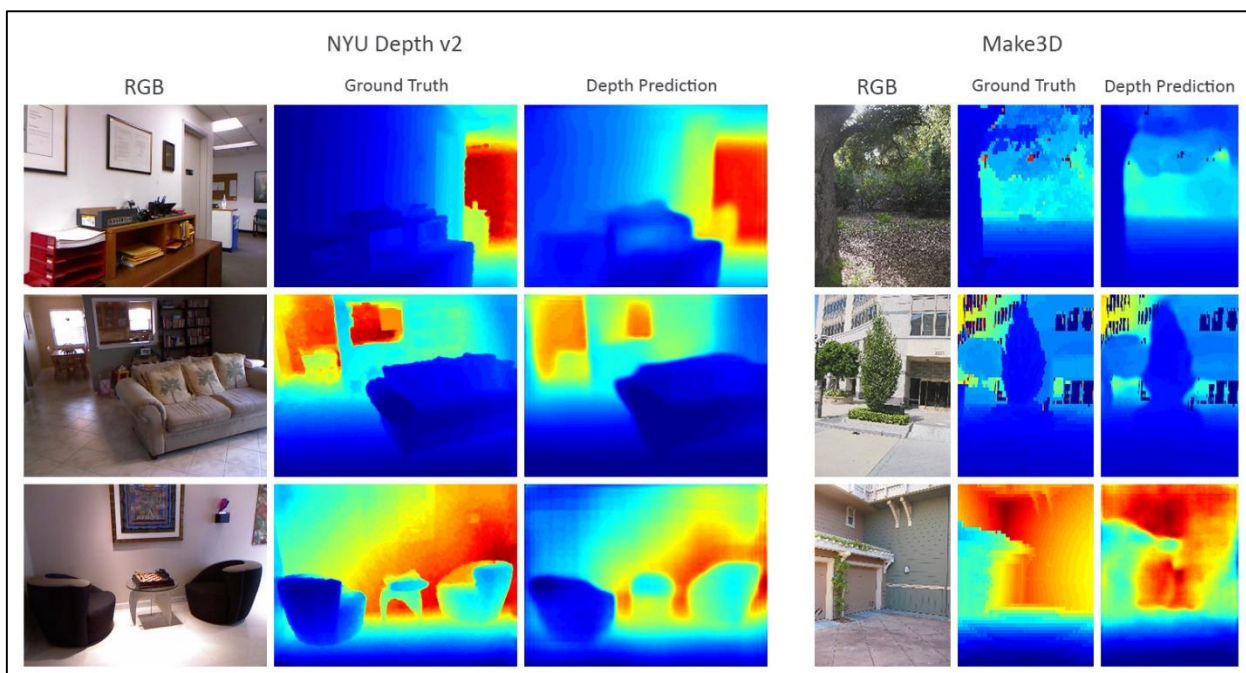


Рисунок 9 – Карта глубины, предсказанная с помощью нейросетей

4.3 Использование учебных данных с аннотацией глубины

В данном методе предполагается, что данные для обучения модели предварительно дополнительно размечены значениями глубины. Он имеет множество преимуществ, таких как: отсутствие необходимости в дополнительных устройствах, высокая мобильность и возможность работы в реальном времени. Но при этом обладает и недостатками: качество напрямую зависит от качество данных для обучения.

4.4 Собственная реализация

4.4.1 Описание модели MiDaS

Для реализации собственного решения была выбрана модель определения глубины MiDaS.

MiDaS (Monocular Depth Estimation via Deep Learning) — это нейросетевая модель, разработанная для оценки глубины сцены из одиночных изображений. Модель была обучена на множестве данных, содержащих изображения с разнообразными сценами, что позволяет ей эффективно определять глубину в различных условиях освещения и композиций.

Принцип ее работы заключается в том, что она использует предобученную на датасете ImageNet сеть для извлечения признаков из входного изображения. Затем, на основе признаков модель предсказывает карту глубины, где каждому пикселю изображения соответствует значение глубины (рисунок 10). Это достигается за счет использования слоев, специально адаптированных для работы с регрессией глубины.



Рисунок 10 – Результаты работы модели MiDaS

4.4.2 Преобразование 2D ключевых точек в 3D координаты

Для трансформации 2D ключевых точек в их 3D координаты используется карта глубины, сгенерированная моделью MiDaS. Каждой ключевой

точке, определённой на изображении, соответствует пиксель на карте глубины. Из этого пикселя извлекается значение глубины, которое представляет собой расстояние от камеры до объекта в данной точке.

На практике, значения глубины могут содержать шумы или неточности из-за ограничений самой модели глубины или из-за визуальных аномалий в исходном изображении. Эти неточности могут привести к ошибкам в определении истинного 3D положения точек. Чтобы уменьшить эти ошибки, применяется процесс коррекции [25]:

- 1) вычисляется евклидово расстояние между соседними ключевыми точками в проекции на плоскость изображения и в пространстве предсказанных глубин;
- 2) если расстояние между точками превышает заданный порог, предполагается, что между точками произошла ошибка в оценке глубины;
- 3) координаты точек корректируются, чтобы минимизировать визуальные искажения, такие как неестественные изгибы или разрывы скелета.

5 Анализ полученных результатов

5.1 Описание эксперимента

В рамках экспериментальной проверки эффективности разработанной модели для распознавания ключевых точек человеческого тела на изображении был реализован комплексный подход, включающий загрузку и обработку изображений, применение предварительно обученной модели и визуализацию результатов. Для достижения максимальной точности и наглядности результатов эксперименты проводились с использованием специализированных библиотек и алгоритмов компьютерного зрения.

Ядро эксперимента составляют PyTorch и его модуль для работы с изображениями torchvision, которые предоставляют инструменты для глубокого обучения и обработки изображений. Это включает в себя загрузку предварительно обученных моделей, операции с тензорами и их обработку.

Для распознавания ключевых точек использовалась модель Keypoint R-CNN, адаптированная для задачи с использованием архитектуры ResNet-50 в качестве основы. Эта модель была обучена на отдельных данных и загружена из сохраненного состояния (checkpoint).

Исходные изображения загружались в формате RGB, после чего преобразовывались в тензоры с помощью функций библиотеки torchvision. Это преобразование включало нормализацию данных и добавление размерности батча для соответствия входным требованиям модели.

Для визуализации результатов распознавания использовались predefined соединения между ключевыми точками (соединения рук, туловища, ног и головы), что позволило наглядно демонстрировать положение и ориентацию человеческого тела на изображении.

Для отрисовки ключевых точек и соединений на изображении использовалась библиотека OpenCV, а результаты визуализировались с помощью Matplotlib.

Для улучшения визуализации трехмерного пространственного распределения ключевых точек человеческого тела, к эксперименту было добавлено

использование технологии визуализации глубины. Это включало применение модели оценки глубины MiDaS. Модель создает карту глубины, которая используется для корректировки и уточнения положения ключевых точек в трехмерном пространстве.

Ключевые точки, распознанные с помощью модели Keypoint R-CNN, в дополнение к двумерным координатам, теперь могли быть обогащены данными о глубине, полученными из карты глубины. Это позволило проводить визуализацию не только в двух измерениях, но и в трехмерной перспективе.

Визуализация трехмерных скелетов проводилась с использованием Matplotlib и его трехмерного модуля *mpl_toolkits.mplot3d*.

5.2 Результаты работы модели определения двумерных ключевых точек

На рисунке 11 можно наблюдать результат работы модели компьютерного зрения, которая выполняет задачу определения позы человека. На изображении изображены четыре человека, каждому из которых модель пыталась присвоить ключевые точки тела и соединить их линиями для визуализации структуры скелета.

Хоть и оценка такого рода результатов является субъективной, кажется, что модель достаточно точно определила ключевые точки всех людей на данном изображении, что говорит об успешности эксперимента.

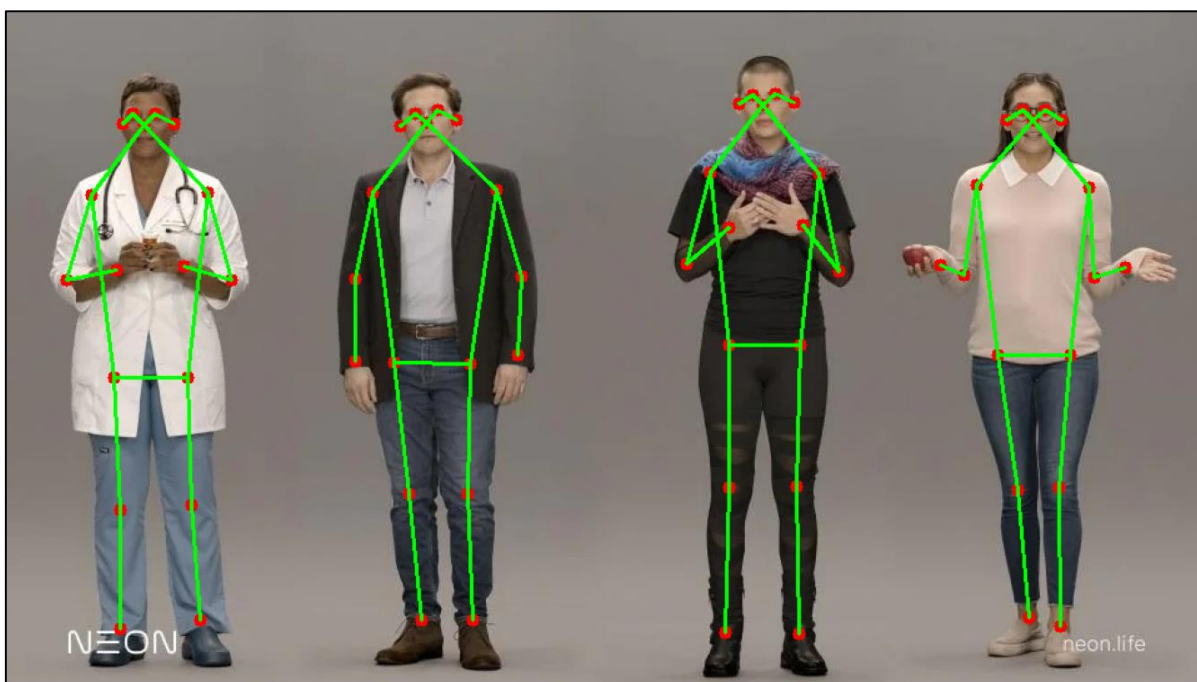


Рисунок 11 – Результат работы собственной модели для определения двумерных ключевых точек

На рисунке 12 изображено изменение скорости обучения (learning rate) по эпохам в процессе обучения нейронной сети. График показывает, как параметр скорости обучения изменяется с течением времени от начала до конца обучения модели.

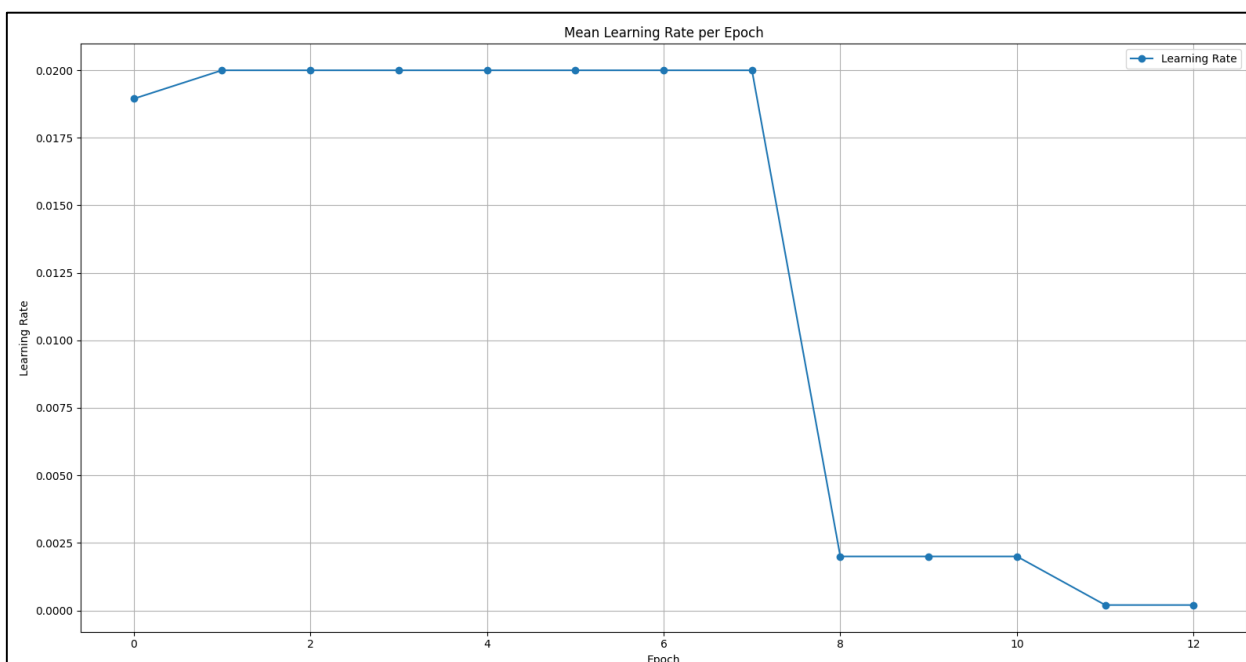


Рисунок 12 – График значений скорости обучения

На начальном этапе обучения, от эпохи 0 до эпохи 4, скорость обучения остается практически неизменной. Это стандартная практика, когда начальные параметры скорости обучения устанавливаются на уровне, который способствует стабильному, но не слишком быстрому улучшению производительности модели.

Наблюдается значительное падение скорости обучения после 5-й эпохи. Это сделано для того, чтобы помочь модели более тонко настраивать свои параметры и избегать переобучения. Особенно это актуально, если начальные эпохи показали достаточное снижение функции потерь и улучшение других метрик точности.

После резкого падения скорость обучения продолжает снижаться и остается на низком уровне в последних эпохах. Это указывает на то, что обучение приближается к завершению, и модель достигает состояния, в котором дальнейшие изменения параметров происходят очень осторожно, чтобы не нарушить уже достигнутую точность.

Снижение скорости обучения во время тренировки модели является общепринятой техникой, которая помогает модели постепенно сходиться к оптимальному решению, уменьшая риск переобучения и позволяя более точно адаптировать веса сети. Этот метод часто используется в сочетании с другими техниками регуляризации и оптимизации процесса обучения.

На рисунке 13 изображено изменение средней потери по ключевым точкам (mean keypoint loss) модели машинного обучения по эпохам обучения. График показывает, как модель оптимизировала свою способность точно определять ключевые точки на объектах в данных.

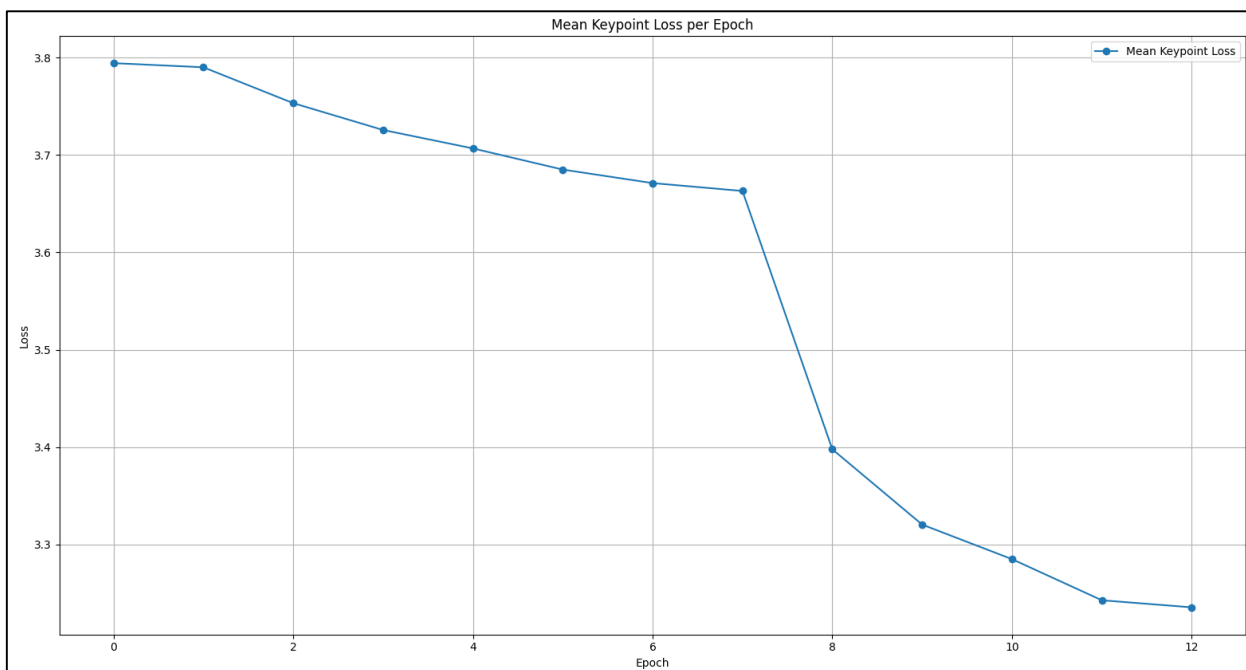


Рисунок 13 – Изменение средней потери по ключевым точкам по эпохам обучения

На начальном этапе обучения средняя потеря ключевых точек постепенно уменьшается. Это указывает на то, что модель начинает адаптироваться к задаче и улучшает свои предсказания благодаря уменьшению разницы между предсказанными и фактическими позициями ключевых точек.

На эпохах 4-8 видно стабильное, но более медленное снижение потерь, что может указывать на постепенное улучшение модели и возможное приближение к пределу своих возможностей на данном наборе данных.

На 8 эпохе можно наблюдать резкое снижение потерь. Это может быть следствием существенного изменения скорости обучения (рисунок 12).

На последних эпохах потери продолжают уменьшаться, хотя и более плавно. Это указывает на то, что модель продолжает оптимизировать свои веса для улучшения производительности, но большинство значительных улучшений уже достигнуто.

Снижение потерь по ключевым точкам на протяжении эпох свидетельствует о эффективности процесса обучения модели. Резкое уменьшение потерь может указывать на успешное применение методик оптимизации обуче-

ния, таких как изменение гиперпараметров или использование адаптивных методов обучения. Стабилизация потерь на более поздних стадиях обучения указывает на то, что модель достигла своего предела точности на доступных данных.

На рисунке 14 изображены изменения средней точности (Average Precision, AP) и средней полноты (Average Recall, AR) для ключевых точек по мере обучения модели.

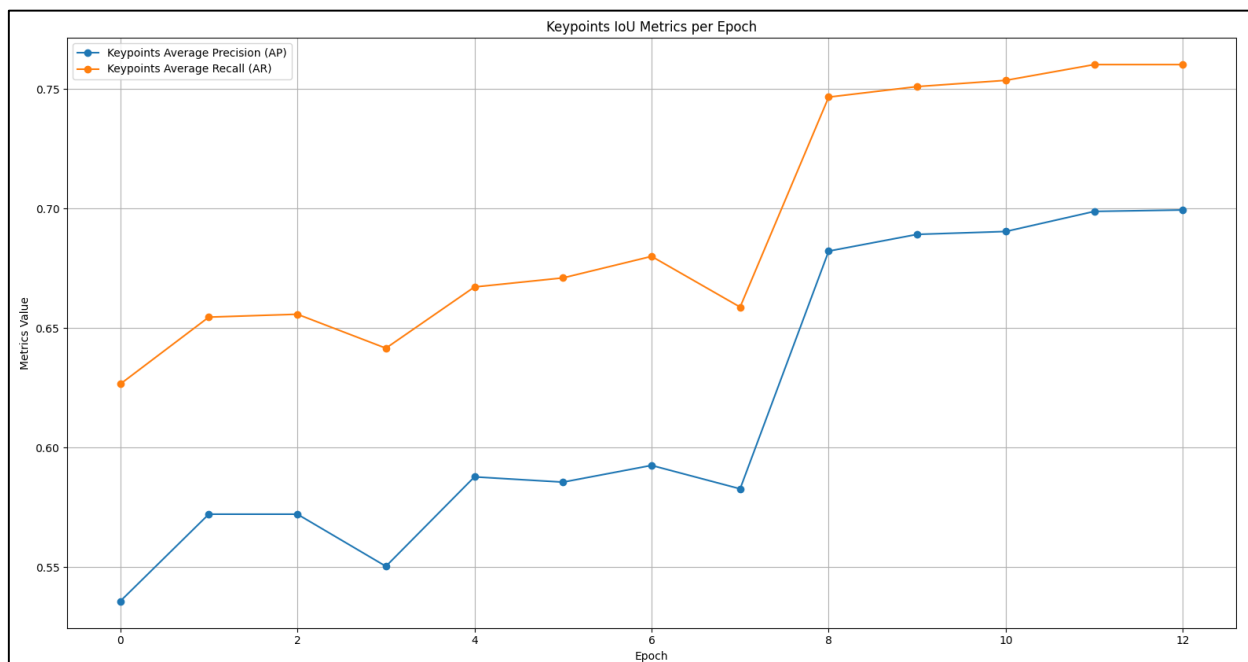


Рисунок 14 – Изменения средней точности и средней полноты для ключевых точек по мере обучения модели

Average Precision (AP) относится к средней точности предсказаний модели. Это мера, которая учитывает и точность, и полноту, где точность — это доля правильно идентифицированных положительных результатов среди всех идентифицированных положительных результатов, а полнота показывает, какая доля реальных положительных результатов была обнаружена моделью.

Average Recall (AR) измеряет способность модели обнаруживать все релевантные случаи в данных. Это доля правильно идентифицированных положительных результатов среди всех реальных положительных случаев.

Начальное значение AP начинается относительно низко, что может указывать на недостаточную способность модели точно локализовать ключевые

точки на ранних этапах обучения. AR начинается выше, чем AP , и демонстрирует более плавные изменения, что характерно для метрики полноты, так как она менее чувствительна к точным локациям и более фокусирована на количестве обнаруженных правильных точек.

Виден рост AP с 6-й по 8-ю эпоху, что коррелирует с резким падением потерь ключевых точек, показанным на предыдущем графике. Это улучшение может быть результатом оптимизаций или корректировки параметров обучения, таких как скорость обучения. Значительное увеличение AR с 6-й по 8-ю эпоху также совпадает с улучшением точности, подтверждая, что модель становится лучше не только в точности, но и в способности обнаруживать больше релевантных ключевых точек.

Стабилизация AP после 8-й эпохи на уровне около 0,70 свидетельствует о достижении моделью некоторого порога эффективности. Постоянство AR на высоком уровне после 8-й эпохи показывает, что модель эффективно обнаруживает большинство ключевых точек.

Эффективность модели в задачах распознавания ключевых точек значительно улучшилась после определенных эпох, что подтверждается увеличением метрик точности и полноты. Управление параметрами обучения, такими как скорость обучения, играет важную роль в достижении оптимальной производительности модели.

5.3 Результаты получения трёхмерных ключевых точек

На рисунке 15 можно увидеть карту глубины, которая получена с помощью MiDaS для дальнейшего использования при построении трехмерного скелета.



Рисунок 15 – Карта глубины, полученная с помощью MiDaS

Карта глубины представляет собой визуализацию, где более темные области соответствуют большей глубине (далее от камеры), а светлые области ближе. Эта карта, полученная из модели глубокого обучения, кажется довольно размытой и нечеткой, что указывает на некоторую неточность в определении точного расстояния до различных частей тела человека. Нечеткость может быть связана с ограничениями самой модели глубокого обучения, такими как недостаточное количество данных для обучения или сложность интерпретации данных сцен с однородными или сложными текстурами.

На рисунке 16 показан трехмерный скелет, созданный на основе карты глубины и определения ключевых точек. Скелет визуализирован с помощью линий, соединяющих ключевые точки, определенные алгоритмом.

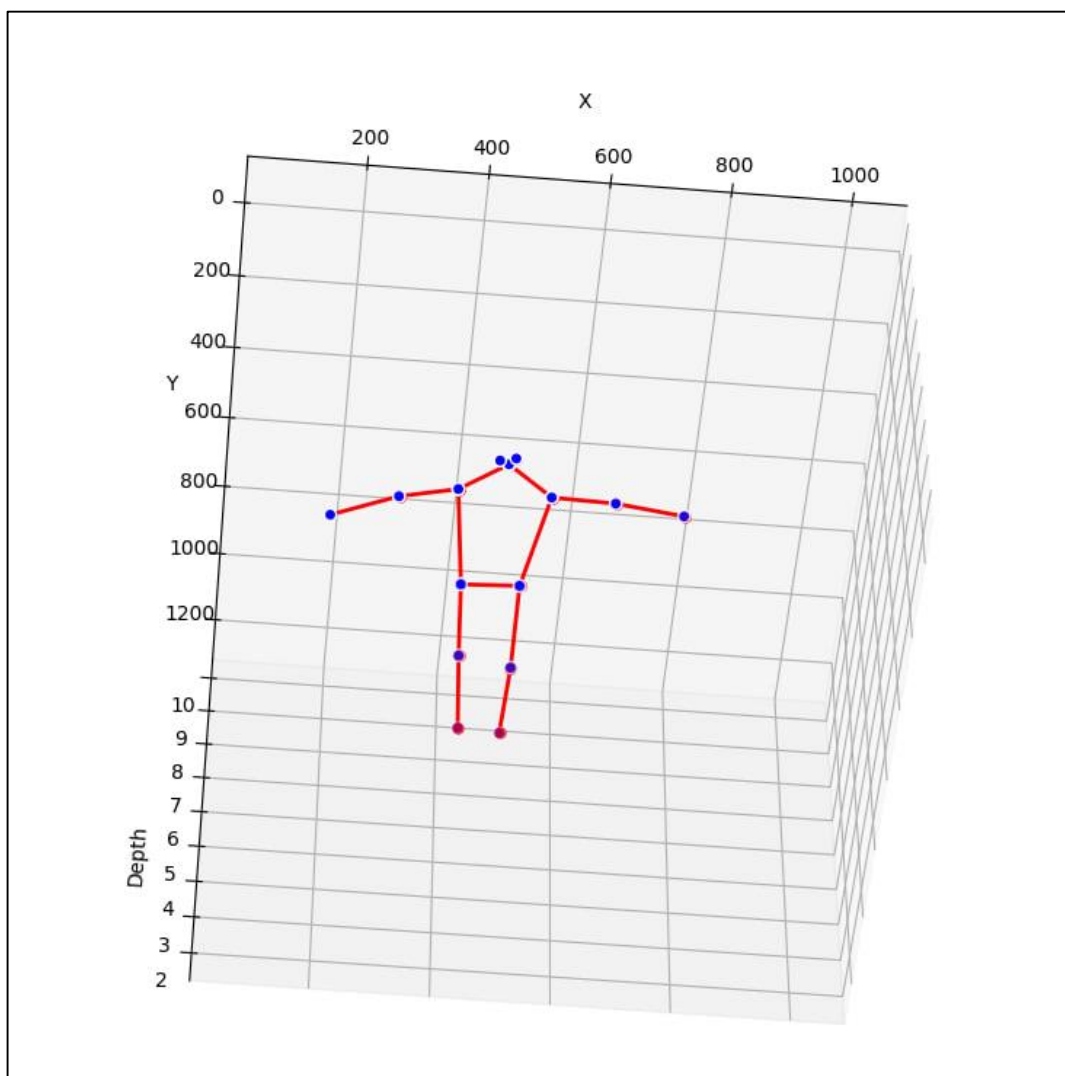


Рисунок 16 – Результат совместной работы собственной модели определения двумерных ключевых точек с моделью MiDaS

Можно заметить, что общая картина трехмерного скелета похожа на правдивую позу из входного изображения. Но, как уже упоминалось, карта глубины демонстрирует значительную нечеткость, что может привести к ошибкам в определении точного положения ключевых точек в пространстве. Это может оказать влияние на точность всей системы построения 3D модели скелета. Ошибки в данных о глубине и определении ключевых точек могут усугубляться при их слиянии. Например, если ключевая точка неправильно

помечена на изображении или если глубина этой точки неверно рассчитана, это приведет к созданию неправильной 3D модели.

Эксперимент подчеркивает значимость точности в каждом из аспектов системы, особенно в генерации карты глубины и определении ключевых точек. Модель для определения 2D ключевых точек демонстрирует хорошие результаты, однако задача генерации карты глубины требует дополнительной работы. Важно отметить, что использование модели глубины MiDaS в данной работе служило лишь в качестве инструмента, а основное внимание было уделено другим аспектам исследования. Существующие проблемы с точностью карты глубины подчеркивают потребность в дальнейших исследованиях и уточнении этой части системы, чтобы достичь более высокой точности в создании трехмерных моделей на основе полученных данных.

ЗАКЛЮЧЕНИЕ

В рамках данной работы было проведено масштабное исследование применимости алгоритмов машинного обучения для захвата и анализа движений человека, особое внимание уделено использованию сверточных нейронных сетей (CNN) и рекуррентных нейронных сетей (RNN) в контексте обработки изображений и видео. Исследование продемонстрировало значительные успехи в области компьютерного зрения, позволяющие точно и эффективно анализировать человеческие движения без использования маркеров.

Основной фокус работы был направлен на разработку и апробацию модели Keypoint R-CNN, адаптированной для задач захвата движения, где ключевую роль играет архитектура ResNet-50 как backbone. Эта модель позволила не только точно локализовать ключевые точки человеческого тела на двумерных изображениях, но и расширить анализ до трехмерного пространства с помощью интеграции алгоритмов оценки глубины изображения.

Среди основных вызовов, с которыми столкнулись в ходе исследования, следует отметить сложности связанные с точностью локализации ключевых точек в условиях низкой контрастности и при частичном перекрытии объектов. Решение данных проблем было частично найдено за счет внедрения корректировки ключевых точек на основе анализа глубины, что позволило значительно повысить точность распознавания движений.

Перспективы развития данной области весьма обширны. Возможность использования алгоритмов для захвата движения без специальных маркеров открывает новые направления в реабилитационной медицине, спортивных науках и развлекательной индустрии. Точный захват движений позволит создавать более сложные и интерактивные системы виртуальной реальности, а также способствовать более глубокому анализу техники спортсменов для предотвращения травм.

Дальнейшие исследования могут быть направлены на улучшение алгоритмов обработки изображений с низким разрешением, а также на разработку новых методов обучения моделей, способных адаптироваться к различным

условиям освещения и фона. Также важным аспектом остается минимизация времени обработки данных в реальном времени для создания систем, способных функционировать в динамичных, непредсказуемых условиях.

Также целесообразным является использование рекуррентных нейронных сетей (RNN), включая LSTM и GRU, для анализа временных последовательностей движения. Эти технологии могут значительно улучшить способность системы распознавать и интерпретировать сложные последовательности движений, что особенно важно в задачах, где необходимо учитывать динамику и предыдущее состояние для предсказания будущих поз.

СПИСОК ИСПОЛЬЗОВАННЫХ ИСТОЧНИКОВ

- 1 Liu W. et al. Large-margin softmax loss for convolutional neural networks // arXiv preprint arXiv:1612.02295. – 2016.
- 2 Романов А. А. Сверточные нейронные сети // Научные исследования: ключевые проблемы III тысячелетия. – 2018. – С. 5-9.
- 3 Yamashita R. et al. Convolutional neural networks: an overview and application in radiology // Insights into imaging. – 2018. – Т. 9. – С. 611-629.
- 4 Осипов В. В. Точечное моделирование операции свертки // Системы. Методы. Технологии. – 2009. – №. 4. – С. 56-63.
- 5 Wu H., Gu X. Max-pooling dropout for regularization of convolutional neural networks // Neural Information Processing: 22nd International Conference, ICONIP 2015, Istanbul, Turkey, November 9-12, 2015, Proceedings, Part I 22. – Springer International Publishing, 2015. – С. 46-54.
- 6 Бабушкина Н. Е., Рачев А. А. Выбор функции активации нейронной сети в зависимости от условий задачи // Инновационные технологии в машиностроении, образовании и экономике. – 2020. – Т. 27. – №. 2. – С. 12-15.
- 7 Javid A. M. et al. A ReLU dense layer to improve the performance of neural networks // ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP). – IEEE, 2021. – С. 2810-2814.
- 8 Basri R. et al. Frequency bias in neural networks for input of non-uniform density // International Conference on Machine Learning. – PMLR, 2020. – С. 685-694.
- 9 Tanhatalab M. R. et al. Deep ran: A scalable data-driven platform to detect anomalies in live cellular network using recurrent convolutional neural network // 2020 IEEE 18th World Symposium on Applied Machine Intelligence and Informatics (SAMI). – IEEE, 2020. – С. 269-274.
- 10 Janai J. et al. Computer vision for autonomous vehicles: Problems, datasets and state of the art // Foundations and Trends® in Computer Graphics and Vision. – 2020. – Т. 12. – №. 1–3. – С. 1-308.

- 11 Lin T. Y. et al. Microsoft coco: Common objects in context // Computer Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13. – Springer International Publishing, 2014. – C. 740-755.
- 12 Deng J. et al. Imagenet: A large-scale hierarchical image database // 2009 IEEE conference on computer vision and pattern recognition. – Ieee, 2009. – C. 248-255.
- 13 Gao S. H. et al. Res2net: A new multi-scale backbone architecture // IEEE transactions on pattern analysis and machine intelligence. – 2019. – T. 43. – №. 2. – C. 652-662.
- 14 Sengupta A. et al. Going deeper in spiking neural networks: VGG and residual architectures // Frontiers in neuroscience. – 2019. – T. 13. – C. 95.
- 15 Mazzini D. Guided upsampling network for real-time semantic segmentation // arXiv preprint arXiv:1807.07466. – 2018.
- 16 Nam S., Lee D. Improvement in object detection using multi-scale RoI pooling and feature pyramid network // Journal of Computing Science and Engineering. – 2022. – T. 16. – №. 1. – C. 14-24.
- 17 Sun Y. et al. Roi pooled correlation filters for visual tracking // Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. – 2019. – C. 5783-5791.
- 18 Wang M. et al. A high-speed and low-complexity architecture for softmax function in deep learning // 2018 IEEE asia pacific conference on circuits and systems (APCCAS). – IEEE, 2018. – C. 223-226.
- 19 Amir I., Koren T., Livni R. SGD generalizes better than GD (and regularization doesn't help) // Conference on Learning Theory. – PMLR, 2021. – C. 63-92.
- 20 Takase T., Oyama S., Kurihara M. Effective neural network training with adaptive learning rate based on training loss // Neural Networks. – 2018. – T. 101. – C. 68-78.

21 Wei L., Zheng C., Hu Y. Oriented Object Detection in Aerial Images Based on the Scaled Smooth L1 Loss Function // Remote Sensing. – 2023. – T. 15. – №. 5. – C. 1350.

22 Van Laarhoven T. L2 regularization versus batch and weight normalization // arXiv preprint arXiv:1706.05350. – 2017.

23 Banks M. S. et al. Stereoscopy and the human visual system // SMPTE motion imaging journal. – 2012. – T. 121. – №. 4. – C. 24-43.

24 Ming Y. et al. Deep learning for monocular depth estimation: A review // Neurocomputing. – 2021. – T. 438. – C. 14-33.

25 Feng M. et al. 2d3d-matchnet: Learning to match keypoints across 2d image and 3d point cloud // 2019 International Conference on Robotics and Automation (ICRA). – IEEE, 2019. – C. 4790-4796.