



**ПРОФЕССИОНАЛЬНОЕ  
ЗАВТРА**



VI Всероссийский сетевой конкурс студенческих проектов с участием студентов с инвалидностью

**МИНИСТЕРСТВО НАУКИ И ВЫСШЕГО ОБРАЗОВАНИЯ  
РОССИЙСКОЙ ФЕДЕРАЦИИ**

**ФГБОУ ВО «Приамурский государственный университет  
имени Шолом-Алейхема»**

**Факультет математики, информационных технологий и техники**

**Направление «Профессиональное завтра в науке»**

**Номинация «Научная статья»**

**«Разработка программы распознавания символов художественного шрифта  
методами машинного обучения»**

**Выполнил:**

**Черкашин Александр Михайлович**

**Руководитель:**

**зав. кафедрой информационных систем,  
математики и правовой информатики  
Баженов Руслан Иванович**

Биробиджан, 2023 г.

Потребность в использовании модели для распознавания дизайнерских символов и художественного шрифта с разными символами возникла при получении текста из изображения. В исследовании предлагается модель машинного обучения, которая выполняет извлекает текст из изображений с разными символами. В наборе данных представлено меню ресторанов и столовых, содержащиеся изображения с готовой разметкой области текста и текст в виде символов. Целью работы является разработка программы и модели по распознаванию и извлечению текста из изображений. Для нахождения текста в изображении использовалась модель ResNet50, а для извлечения текста применялась модель DenseNet. Предлагаемый метод обеспечивает высокие показатели точности при распознавании и нахождении текста в изображениях. Таким образом, предложенная модель извлечения символов художественных шрифтов с использованием машинного обучения является весьма перспективной и имеет большое значение для разработки новых моделей распознавания в области компьютерной зрения. Масштабы результатов широки, и полученные результаты могут быть применены ко многим практическим приложениям, включая графический дизайн, типографику и цифровые медиа.

Ключевые слова: Python, Pytorch, Tensorflow, OCR, Recognizer, Detector, ResNet, DenseNet.

## 1. Обзор исследований

И. Цуляк и др. описывают различные компоненты OpenCV и его применения в обработке изображений и видео, компьютерном зрении и машинном обучении. Исследовательская работа также включает практические примеры и фрагменты кода, которые подчеркивают использование функций и функций OpenCV для задач обработки изображений и видео. В целом авторы стремятся дать новичкам в области компьютерного зрения базовое представление о библиотеке OpenCV и ее потенциальных приложениях [1].

К. Пулли и др. исследуют важность эффективного и точного компьютерного зрения в реальном времени в различных областях, таких как робототехника, видеонаблюдение, анализ видео и медицинская визуализация. Затем они дают обзор различных методов и инструментов, используемых в компьютерном зрении в реальном времени, и того, как библиотека OpenCV может использоваться в этих приложениях [2].

Исследование, проведенное Г. Хие, В. Лу посвящено обнаружению границ на изображениях с использованием технологии OpenCV. Авторы представляют важность обнаружения краев в различных приложениях обработки изображений. Затем объясняют концепцию краев изображения и общие методы, используемые для их обнаружения [3].

Г. Агам предоставляет пошаговое руководство по установке OpenCV, настройке среды разработки и началу программирования на C++. Автор обсуждает фундаментальные концепции обработки изображений, обнаружения объектов и извлечения признаков, а также показывает, как реализовать эти методы с помощью различных функций и модулей OpenCV. Кроме того, учебник включает примеры реальных приложений компьютерного зрения, таких как распознавание лиц, обнаружение движения и отслеживание транспортных средств. В целом автор стремится предоставить исчерпывающее введение в программирование OpenCV,

доступное для начинающих, а также полезное для опытных разработчиков [4].

Л. Яиао и др. предоставляет обзор различных инструментов, включенных в DavarOCR, включая инструменты распознавания текста, инструменты распознавания таблиц и инструменты распознавания изображений. Инструментарий разработан таким образом, чтобы быть простым в использовании и адаптируемым к широкому спектру вариантов использования. Авторы описывают ряд тематических исследований и экспериментов, демонстрирующих эффективность DavarOCR. Эти исследования показывают, что DavarOCR способен точно распознавать и понимать широкий спектр документов и может работать с различными языками. В целом цель авторов состоит в том, чтобы предоставить мощный и гибкий инструмент, который может помочь исследователям и практикам лучше понять и проанализировать широкий спектр документов [5].

В исследовании М. Висинтин и А. Девер. по оптическому распознаванию символов (OCR) представлен обзор технологии, используемой для автоматического распознавания печатного или рукописного текста в цифровых изображениях или отсканированных документах. Авторы обсуждают различные методы и алгоритмы, которые используются для извлечения, анализа и интерпретации текста, включая предварительную обработку, сегментацию, извлечение признаков и машинное обучение. Они также изучают применение оптического распознавания символов в различных областях, таких как оцифровка документов, индексация и поиск текста, а также автоматический перевод [6].

Н. Лу и др. предлагают многоаспектную нелокальную сеть (Master) для оптимального распознавания текста сцены. Они утверждают, что существующие методы распознавания текста сцены ограничены, поскольку они рассматривают задачу распознавания как простую проблему классификации изображений, которая не учитывает сложные вариации и искажения, возникающие в естественных сценах. Чтобы устранить это

ограничение, авторы предлагают новую структуру, включающую нелокальные нейронные сети, которые позволяют сети фиксировать долгосрочные зависимости между различными частями входного изображения. Система обучается на большом наборе данных текстовых изображений с использованием комбинации контролируемых и неконтролируемых методов обучения для повышения точности и обобщения. [7].

В. Ю и др. представляет метод извлечения ключевой информации из документов с использованием графовых сверточных сетей (GLCN). Авторы предлагают полуконтролируемую структуру, которая сочетает в себе явное распространение сигнала с выделением признаков для повышения точности извлечения информации. Предлагаемый метод оценивается на нескольких эталонных наборах данных, и результаты показывают, что он превосходит несколько современных методов с точки зрения точности, полноты и оценки F1. В документе также представлен подробный анализ предлагаемого метода, включая изучение влияния различных гиперпараметров на производительность. В целом исследование показывает, что предложенный ими метод является эффективным подходом к сложной проблеме извлечения ключевой информации из документов [8].

Авторы исследования П. Х. Лин и Д. Ван стремятся повысить точность и эффективность систем OCR, разработав модель, позволяющую редактировать шаблоны, то есть процесс манипулирования шаблонами OCR для исправления ошибок распознавания текста. Система включает в себя редактор шаблонов, который позволяет пользователям вручную редактировать шаблоны OCR, и алгоритм проецирования слов, который предсказывает следующее слово в предложении, повышая точность системы распознавания. В целом исследование направлено на повышение точности и эффективности систем OCR, делая их более полезными и надежными инструментами для распознавания текста [9].

Авторы исследования З. Куанг и др. стремятся создать гибкую и настраиваемую структуру, которую можно применять к целому ряду приложений, от оцифровки документов до поиска изображений. Исследование включает в себя разработку ряда моделей глубокого обучения, которые могут точно и эффективно обнаруживать, и распознавать текст на изображениях с использованием различных методов, включая сверточные нейронные сети, механизмы внимания и пространственные преобразователи. В целом, исследование направлено на совершенствование современных технологий распознавания текста и предоставление исследователям и разработчикам комплексного набора инструментов для работы с текстом в изображениях [10].

О. Зайен и др. приняли участие в конкурсе ICPR2020 и разработали подход, основанный на глубоком обучении, для обнаружения текстовых областей в новостных видеокдрах на арабском языке. Они применили методы предварительной и последующей обработки для повышения производительности модели. Кроме того, авторы предложили систему распознавания для распознавания и сегментации текстовых областей для достижения оптимальных результатов. Основная цель исследования — улучшить систему обнаружения и распознавания арабского текста в наборе данных новостных видеороликов путем разработки надежного подхода, основанного на глубоком обучении [11].

Г. Яуме, Х. К. Екенел, Й. П. Тхиран представляют новый набор данных под названием FUNSD, что означает «Понимание форм в зашумленных отсканированных документах». Этот набор данных предназначен для продвижения исследований в области понимания формы и извлечения информации из отсканированных документов. В исследовании приводится подробное описание набора данных, включая его содержание и методологию аннотирования. Авторы также представляют результаты тестов, основанные на различных подходах к машинному обучению с использованием набора данных FUNSD. Кроме того, они дают представление о проблемах,

связанных с набором данных, и потенциальных будущих направлениях исследований в этой области. В целом, исследование направлено на то, чтобы предоставить ценный ресурс для исследователей, заинтересованных в развитии области понимания формы и извлечения информации из отсканированных документов [12].

Авторы П. Батра и др. стремятся повысить точность и скорость систем ALPR с помощью алгоритма You Only Look Once (YOLO). Этот алгоритм был модифицирован и обучен на большом наборе данных номерных знаков для распознавания номерных знаков в режиме реального времени с высокой точностью. Исследователи провели эксперименты, чтобы продемонстрировать эффективность и действенность предложенной ими системы ALPR. В целом, представляют новый подход к системам ALPR, который может быть полезен в различных приложениях, таких как управление дорожным движением, управление парковками и системами безопасности [13].

М. Памнани рассказывает о студенческом проекте по разработке приложения для Android под названием DOT-HAZMAT. Приложение предназначено для использования технологий машинного обучения и компьютерного зрения для поиска и классификации табличек HAZMAT во время спасательных операций. Автор описывает приложение как инновационное решение для повышения безопасности спасателей во время аварий с опасными материалами. Проект подробно объясняется, демонстрируя технические аспекты технологии и ее потенциальное использование в будущем. В целом, автор представляет обзор студенческого проекта и его вклада в улучшение мер безопасности во время инцидентов HAZMAT [14].

А. Весалаинен фокусируется на разработке и оценке различных методов сегментации изображений для детального анализа макета документа OCR. Сначала автор проводит литературный обзор предыдущих исследований по анализу макета документа и сегментации изображений, а

также определяют ограничения и проблемы существующих методов. Затем предлагает и сравнивает различные подходы к сегментации изображений документа на области текста, изображения и фона с использованием таких методов, как адаптивное пороговое значение, обнаружение границ и анализ связанных компонентов. Автор оценивает точность и эффективность каждого метода, используя большой набор данных документов OCR, и обсуждает свои выводы и рекомендации для будущих исследований в этой области. В целом исследование направлено на то, чтобы внести свой вклад в разработку более точных и универсальных систем оптического распознавания символов, способных обрабатывать сложные макеты документов [15].

Авторы работы К. Олейничак и М. Шульц фокусируются на выявлении и анализе проблемы обнаружения текста в системах оптического распознавания символов (OCR). Авторы утверждают, что большинство современных систем оптического распознавания текста в первую очередь ориентированы на распознавание текста из отсканированных документов, но не обязательно на обнаружение текста в изображениях документов. В результате эти системы могут пропускать или неправильно распознавать текст на изображениях с низким разрешением или сложным фоном, что приводит к снижению точности оптического распознавания символов. Для решения этой проблемы авторы предлагают новый подход, сочетающий методы обнаружения текста с системами OCR. Они проводят эксперименты, чтобы сравнить точность этого подхода с существующими системами OCR и демонстрируют, что предлагаемый ими подход повышает точность распознавания для изображений со сложным фоном или разрешением низкого качества. В исследовании подчеркивается важность включения распознавания текста в системы OCR для более точного распознавания текста на изображениях документов [16].



## 2. Исходные данные

Исходные данные представлены на сайте конкурса Huawei digix global ai challenge 2021 в категории Menu identification leveraging few-shot learning models [17].

Исходный датасет содержится в виде двух файлов архива 2021\_5\_data.zip и 2021\_finals\_data\_5.zip. В файле архива включаются каталоги train\_image\_common (3561 изображения), test\_image (500 изображения) и train\_image\_special (93 изображения), в которых содержится изображением с наложенным текстом, а файлы train\_label\_common.json и train\_label\_special.json содержат указатели с выделенными изображениями, с размеченными областями и текстом. Всего содержится 4154 изображения.

Архивы 2021\_5\_data.zip и 2021\_finals\_data\_5.zip отличается только в каталоге test\_image тестовые данные (Рис 1).

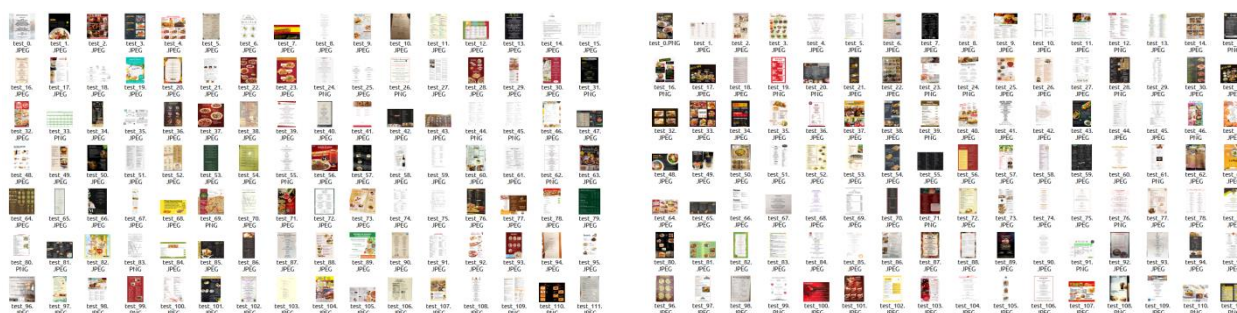


Рисунок 1. Отличие набор данных, слева набор данных из архива 2021\_5\_data.zip, а справа 2021\_finals\_data\_5.zip, в каталоге test\_image

Таблица 1. Архив 2021\_5\_data.zip и 2021\_finals\_data\_5.zip

Название	Описание
train_image_common	Тренировочные данные для машинного обучения
test_image	Тестовые данные для оценки качества модели
train_image_special	Специальные тренировочные данные для оценки качества модели
train_label_common.json	Файл выделенных тренировочных данных с изображениями текста

train\_label\_special.json

Файл выделенных специальных  
тренировочных данных с  
изображениями текста

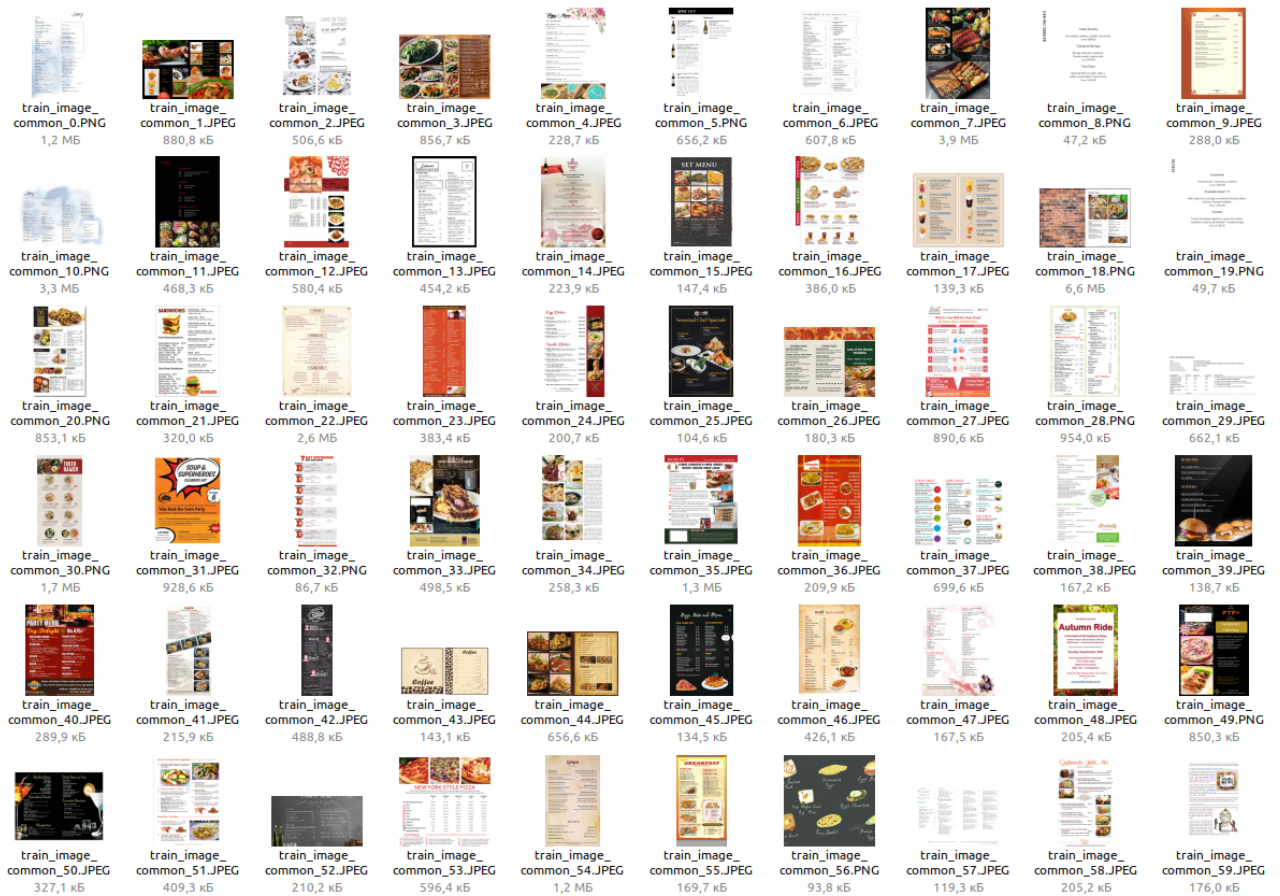


Рисунок 2. Исходные данные в каталоге train\_image\_common

Файл train\_label\_common.json содержит ассоциативный массив, ключ — название файла изображения из каталога train\_image\_common, значение структура (табл 2).

Таблица 2 Файлы train\_label\_common.json и train\_label\_special.json

Название	Описание
label	Текст
points	Вершины с координатами x, y с плавающей точкой

```
IPython: Milfslm/2021_5_data
sasha@sasha-ecoblue: /home/ram/ram/Milfslm/2021_5_data$ ipython3
Python 3.10.4 (main, Apr 2 2022, 09:04:19) [GCC 11.2.0]
Type 'copyright', 'credits' or 'license' for more information
IPython 8.3.0 -- An enhanced Interactive Python. Type '?' for help.

In [1]: import json

In [2]: f = open("train_label_common.json", "r")

In [3]: js_a = json.load(f)

In [4]:
```

Рисунок 3. Чтение файла train\_label\_common.json

```
IPython: Milfslm/2021_5_data

In [3]: js_a = json.load(f)

In [4]: js_a["train_image_common_1.JPEG"]
Out[4]:
[{'label': 'FOR YOU ?',
  'points': [[242, 746.279296875],
             [323.2000000000007, 746.279296875],
             [323.2000000000007, 767.800000000011],
             [242, 767.800000000011]]},
 {'label': 'Chef Spicy Sliced Beef',
  'points': [[429, 740.2998046875],
             [587, 740.2998046875],
             [587, 758.700000000007],
             [429, 758.700000000007]]},
 {'label': 'KETCHUP',
  'points': [[249, 728.0263671875],
             [317, 728.0263671875],
             [317, 746.300000000011],
             [249, 746.300000000011]]},
 {'label': 'CHILLI OR',
  'points': [[246.800000000011, 707.7998046875],
             [319, 707.7998046875],
             [319, 725.9998046874989],
             [246.800000000011, 725.9998046874989]]},
```

Рисунок 4. Вывод структуры из файла train\_label\_common.json. Выбранный «train\_image\_common\_1.JPEG»



Рисунок 5. Файл train\_image\_common\_1.JPEG. Расположенный  
2021\_5\_data/train\_image\_common

Структура файлов train\_label\_common.json и train\_label\_special.json:

- Имя файла изображение: Словарь
  - Список
    - label: Строка
    - points: Список
      - Список: по координаты X и Y.

### 3. Структура программы

Для работы получено базовое решение, исходный код программы 2021-DIGIX-BASELINE в каталоге baseline-game5 [18].

Таблица 3. Файловая структура вывода в программе содержится

Адрес	Описание
output/checkpoints	Содержатся сохраненные и обученные модели
output/checkpoints/detector/checkpoint.pth.tar	Файл, обученный модель (Detector) для поиска текста изображения
output/checkpoints/recognizer	Содержатся серии сохраненные и обученные модели для извлечение текста из изображения (Recognizer)
output/detector_test_output/menu	Содержатся серии изображения выведено выделенный области (нарисованные квадратики) нахождения текста. Как результат

	обученный модель по поискам текста изображения (Detector).
output/recognizer/labels	Содержатся файлы формата json, в структуре файла содержится область нахождения (вершины) текста и текст в результате полученная обученная модель по извлечению текста из изображения (Recognizer).
output/recognizer_images	Содержатся серии изображения, вырезанные из разметки в файла train_label_common.json и train_label_special.json содержащие вершины. Изображение содержит только текст.
output/recognizer_log/train	Содержатся серии журнала обучения модели по извлечению текста из изображения (Recognizer).
output/recognizer_txts/train.txt	Файл содержит список в виде строки, имя изображение (из серии output/recognizer_images) и текст значение.
output/recognizer_txts/real_train.txt	Файл содержит список в виде строки, имя изображение (из серии output/recognizer_images) и текст значение в виде коды символы.
output/test_null.json	Файл в формате json, в структуре файла содержится область нахождения (вершины) текста и пустой текст.

Программа написана на языке Python, которая использует библиотеки PyTorch, OpenCV, Keras, Tensorflow, для ускорения обработки изображения программная часть была написано на языке C++ с использованием библиотеки OpenCV. Для ускорения работы по предсказанию применяется библиотека multiprocessing.

### 3.1. Запуск detector/train.py

В этой программе выполняется поиск текст изображений (Detector).

Использовали модель ResNet50 и оптимизацию SGD для обучения по нахождению текста в изображения, для оценки обучения применяется метрика Dice (Рис. 2) [19].

$$\text{Dice} = \frac{2 * TP}{2 * TP + FP + FN}$$

Рисунок 6. Формула метрики Dice

### Параметры оптимизация

```
SGD (
Parameter Group 0
  dampening: 0
  lr: 0.0001
  maximize: False
  momentum: 0.99
  nesterov: False
  weight_decay: 0.0005
)
```

### Модель ResNet50

```
OCR_DETECTOR(
(backbone): ResNet(
  (conv1): Conv2d(3, 64, kernel_size=(7, 7), stride=(2, 2), padding=(3, 3), bias=False)
  (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  (relu1): ReLU(inplace=True)
  (maxpool): MaxPool2d(kernel_size=3, stride=2, padding=1, dilation=1, ceil_mode=False)
  (layer1): Sequential(
    (0): Bottleneck(
      (conv1): Conv2d(64, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
      (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
      (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
      (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      (relu): ReLU(inplace=True)
      (downsample): Sequential(
        (0): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
        (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
      )
    )
  )
  (1): Bottleneck(
    (conv1): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
  (2): Bottleneck(
    (conv1): Conv2d(256, 64, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(64, 64, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
  )
)
```

```

(conv3): Conv2d(64, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
)
(layer2): Sequential(
  (0): Bottleneck(
    (conv1): Conv2d(256, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): Conv2d(256, 512, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
  (1): Bottleneck(
    (conv1): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
  (2): Bottleneck(
    (conv1): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
  (3): Bottleneck(
    (conv1): Conv2d(512, 128, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(128, 128, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(128, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
  )
)
)
(layer3): Sequential(
  (0): Bottleneck(
    (conv1): Conv2d(512, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
    (bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    (relu): ReLU(inplace=True)
    (downsample): Sequential(
      (0): Conv2d(512, 1024, kernel_size=(1, 1), stride=(2, 2), bias=False)
      (1): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
    )
  )
  (1): Bottleneck(

```

```

(conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(2): Bottleneck(
(conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(3): Bottleneck(
(conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(4): Bottleneck(
(conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(5): Bottleneck(
(conv1): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(256, 1024, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(1024, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
)
(layer4): Sequential(
(0): Bottleneck(
(conv1): Conv2d(1024, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
(downsample): Sequential(
(0): Conv2d(1024, 2048, kernel_size=(1, 1), stride=(2, 2), bias=False)
(1): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
)
)
(1): Bottleneck(
(conv1): Conv2d(2048, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)

```



```

(conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
(2): Bottleneck(
(conv1): Conv2d(2048, 512, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv2): Conv2d(512, 512, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
(bn2): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(conv3): Conv2d(512, 2048, kernel_size=(1, 1), stride=(1, 1), bias=False)
(bn3): BatchNorm2d(2048, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu): ReLU(inplace=True)
)
)
)
)
(neck): FPN(
(reduce_layers): Sequential(
(0): Sequential(
(0): Conv2d(256, 256, kernel_size=(1, 1), stride=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
(1): Sequential(
(0): Conv2d(512, 256, kernel_size=(1, 1), stride=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
(2): Sequential(
(0): Conv2d(1024, 256, kernel_size=(1, 1), stride=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
(3): Sequential(
(0): Conv2d(2048, 256, kernel_size=(1, 1), stride=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
)
(smooth_layers): Sequential(
(0): Sequential(
(0): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
(1): Sequential(
(0): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
(2): Sequential(
(0): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
(1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(2): ReLU(inplace=True)
)
)
)
)
(head): SIMPLE_DILATE_HEAD(
(conv1): Conv2d(1024, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
(bn1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True, track_running_stats=True)
(relu1): ReLU(inplace=True)

```

```
(conv2): Conv2d(256, 7, kernel_size=(1, 1), stride=(1, 1))
)
```

Аугментация данных выполняется перед каждый шаг обучения.

Этапы аугментация изображения:

1. Используем трех-канальные RGB.
2. Изменяем случайные размер изображение по масштабу от 0.5 до 3.0.
3. Создаем два изображения, текст и маска. Текст — черный, маска (область изображения) — белый фон.
4. Рисуем выделенный область для текста (points), если текст пустой или # (решетка) то накладываем черную область в изображение маска.
5. Создаем несколько изображений маски с названием ядро, в данном случае 7 слоев (Рис 5).
6. Обрабатываем изображение, уменьшаем размер до 640 пикселей, случайно поворачиваем и обрезаем изображение (Рис 6).
7. Изменяем случайно яркость (в значение 32 / 255) и насыщенность (в значение 0.5).
8. Преобразуем в тензор.
9. Нормализуем изображения.

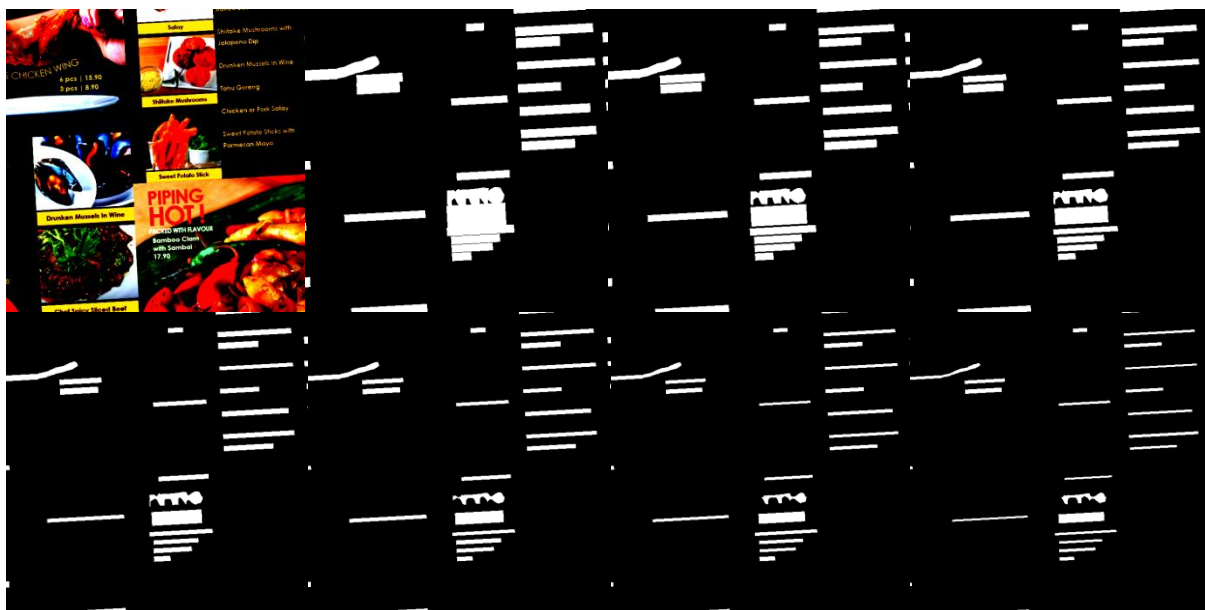


Рисунок 7. Серия изображение вырезано размером 640x640. Вверхни  
слева стороны изображение

Левый и верхний рисунок - обработанное исходное изображения (готов к обучению), следом черно-белый рисунок — маска под текст, выделенное белое (Маска) геометрический вершины уменьшается в центр, начиная от 0.9 заканчивая до 0.3 (Рис 7). Область обучения (Маска обучение), цвет области устанавливается, если исходные изображение обработано и пересекается за край, белый — распознаваемый область, черный — не распознаваемая область, то есть исходное изображение не лежит в этой области (Рис 8).



Рисунок 8. Левоы исходный изображений, справа область распознаваемый изображений (маска распознавания)

```
sasha@sasha-ecoblue: /home/ram/2021-DIGIX-BASELINE/baseline-game5/dete...
sasha@sasha-ecoblue: /home/ram/2021-DIGIX-BASELINE/baseline-game5/detector$ python3 ./train.py
****model info****:
=> 1.backbone: resnet50
=> 2.neck: fpn
=> 3.head: simple_dilate_head
****TRAIN DATASET****:
1. menu_train: ***sum_samples: 3561
=> Training Mode: FROM-SCRATCH
15
5
Using step LR schedule with warm-up epochs of 0!
****start training****:
=> 1.Starting Epoch: 0
=> 2.Total Epoch: 15
/mnt/ram/2021-DIGIX-BASELINE/baseline-game5/detector/dataset/train_dataset.py:95: UserWarning: Creating a
tensor from a list of numpy.ndarrays is extremely slow. Please consider converting the list to a single
numpy.ndarray with numpy.array() before converting to a tensor. (Triggered internally at ../torch/csrc/
utils/tensor_new.cpp:210.)
text_mask = torch.tensor(target[0: kernel_num]).ge(0.5).float()

=>Epoch 0, learning rate = 0.0001
batch: 1/3561, loss: 0.3292. iou-t: 0.328. iou-k: 0.457. accuracy: 0.433
/mnt/ram/2021-DIGIX-BASELINE/baseline-game5/detector/dataset/train_dataset.py:179:
VisibleDeprecationWarning: Creating an ndarray from ragged nested sequences (which is a list-or-tuple of
lists-or-tuples-or ndarrays with different lengths or shapes) is deprecated. If you meant to do this, you
```

Рисунок 9. Обучения модели

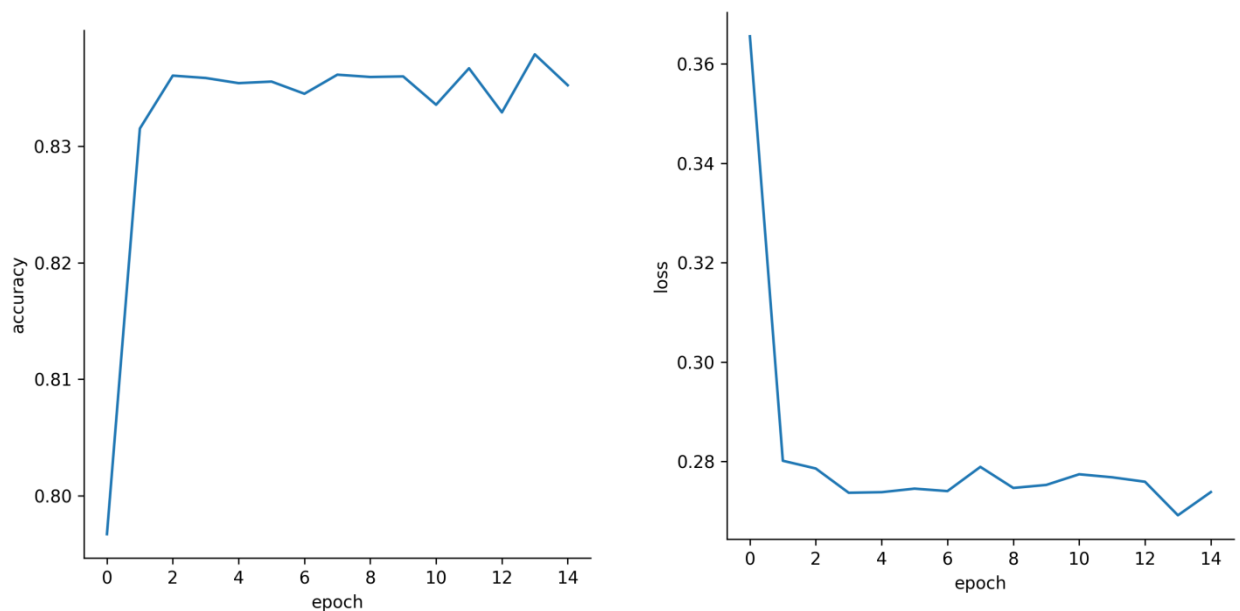
Таблица 4. Вывода данные метрики, batch 3561.

Epoch (Эпоха)	Loss (Dice)	iou-t	iou-k	accuracy
0	0.365571	0.634908	0.623926	0.796715
1	0.280158	0.714952	0.657267	0.831510
2	0.278601	0.719787	0.660720	0.836051
3	0.273721	0.723036	0.661050	0.835848
4	0.273836	0.721631	0.659525	0.835405

Epoch (Эпоха)	Loss (Dice)	iou-t	iou-k	accuracy
5	0.274550	0.721451	0.659668	0.835540
6	0.274054	0.720512	0.661375	0.834495
7	0.278926	0.721446	0.662046	0.836130
8	0.274689	0.722450	0.660161	0.835934
9	0.275296	0.721505	0.661585	0.835982
10	0.277442	0.718966	0.659482	0.833552
11	0.276836	0.721942	0.662015	0.836680
12	0.275936	0.719033	0.658878	0.832896
13	0.269190	0.724763	0.663797	0.837875
14	0.273860	0.720973	0.661448	0.835225

Iou-t (Intersection over Union - text) — оценка сходства текстовое область изображение (маска текста) в модели обучение (табл. 4).

Iou-k (Intersection over Union - kernel) — оценка сходства ядро область изображение (маска область распознавания) в модели обучения (табл. 4).



Accuracy — оценка точности модели для маски текста (табл. 4).

Рисунок 10. На графике представлены левой accuracy, Справа Loss. По оси Ох - epoch (Эпоха)

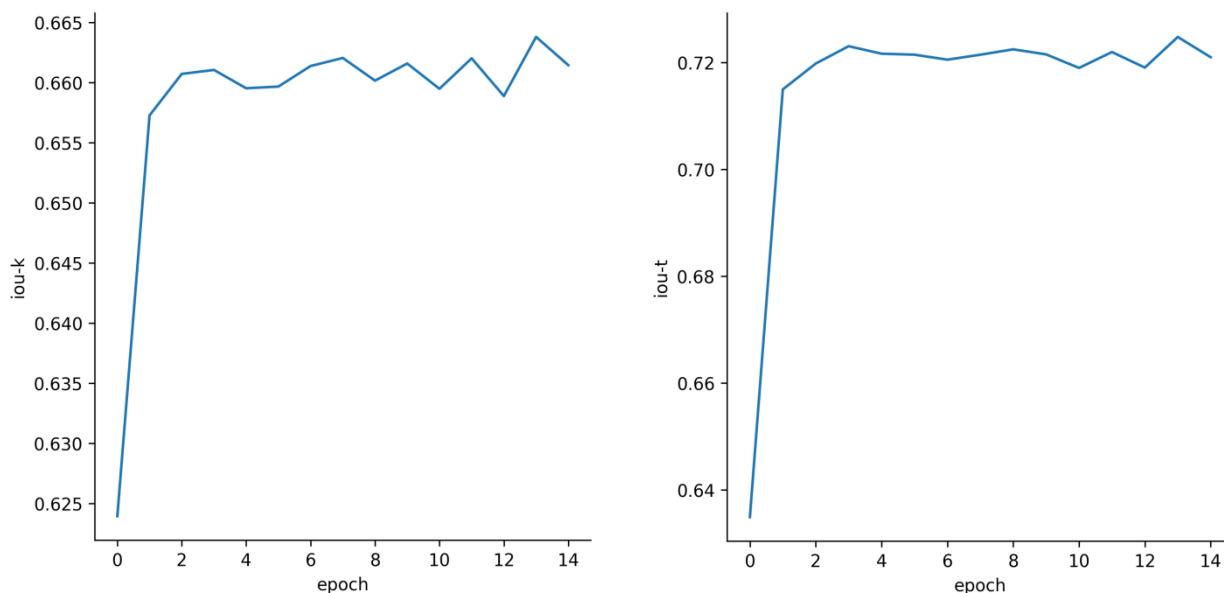


Рисунок 11. На графике представлены левой iou-k, Справа iou-t. По оси Ох - epoch (Эпоха)

Время выполнение программы (учитывается обучение и мониторинг метрики, вывод в консоль и аугментация данных) равно 32136.69 секунд (8 часов, 55.61 минут, 36.69 секунда). Время считается с начало и до конца цикла обучения в каждый эпохи обучения.

Средняя время шаг обучения 0,602 секунда.

Время выполнение обучения (не учитывается мониторинг метрики, вывод в консоль и аугментация данных) равно 15293.21 секунд (4 часов, 14 минут, 53,21 секунда). Время считается с начало и до конца шага обучения в каждый цикл обучения.

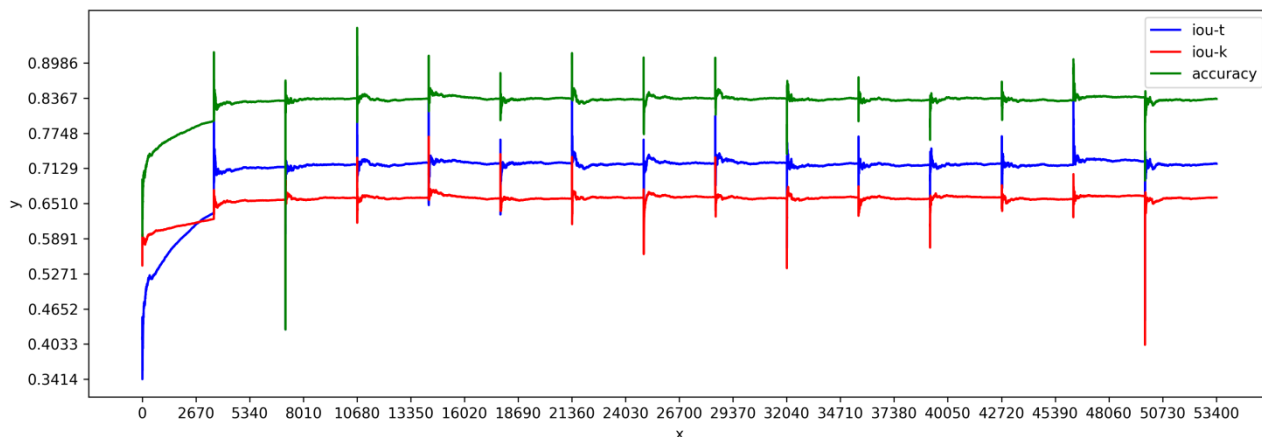


Рисунок 12. На графике представлены цикл обучение за все эпохи

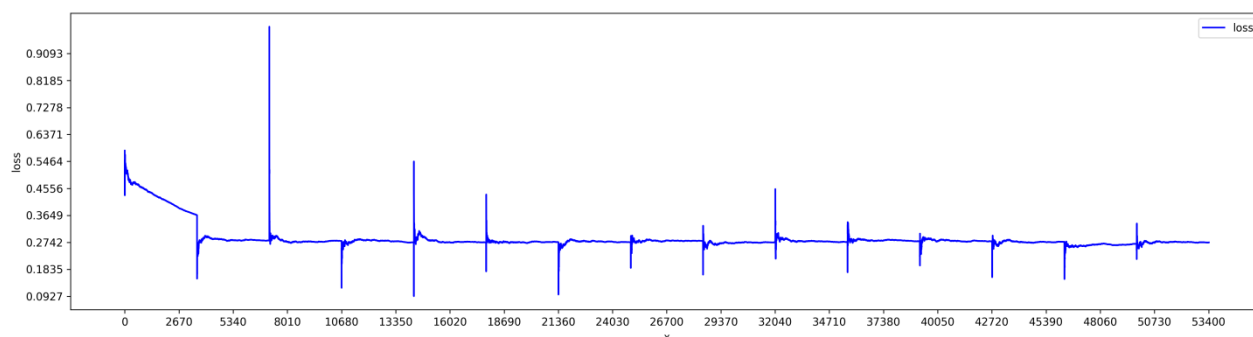


Рисунок 13. На графике представлены цикл обучение за все эпохи с функцией потерь (Dice)

На графике представлены у это метрики accuracy, iou-t, iou-k, loss, а x это цикл обучение за все эпохи (Рис 12, 13). Максимальный цикл обучение за все эпохи равен 53 415.

### 3.2. Подготовка данных к обучению recognizer

В этой программе подготавливается исходные данные для обучения по извлечению текста из изображений.

1. Чтение файла `train_label_common.json` и `train_label_special.json` в формате json.
2. Чтение изображения `train_image_common` и `train_image_special` указанный в `train_label_common.json` и `train_label_special.json`.
3. Получаем points и сортируем вершины по возрастанию.





\*  
+  
,  
-  
.  
/

5. Кодруем символы в число.

6. Выводим текстовый файл real\_train.txt.

```
0.jpg 34 41 44 44 37 35 33 50 52 0 51 33 44 45 47 46 0 34 50 53 52 0 50 47 51 37
1.jpg 41 48 33
2.jpg 40 41 52 33 35 40 41 46 47 0 46 37 51 52 0 36 33 41 0 36 33 41
3.jpg 43 41 50 41 46 0 36 50 33 53 39 40 52
4.jpg 50 72 66 71 12 0 69 81 64 70 81 64 77 83 12 0 84 76 64 76 72
5.jpg 39
6.jpg 51 71 72 72 83 64 74 68 0 45 84 82 71 81 78 78 76 82 0 86 72 83 71
7.jpg 34 33 50 43 0 35 33 38 37 0 38 33 45 47 53 51 0 35 40 41 35 43 37 46 0 55 41 46 39
8.jpg 66 81 64 89 88
```

Время выполнение программы (от запуска до завершения) равно 1 минуте и 15,565 секунд.

В результате вывода набор вырезанных изображений содержит 30 164 изображений (Рис 15).



Рисунок 16. Вырезанный изображений

### 3.3. Запуск recognizer/train.py

В программе выполняется обучение модели по распознаванию текста в изображениях (Recognizer). Используется библиотека Keras.

Параметры обучения:

- Размер пакетов `batch_size = 128`.
- Максимальная длина строки текста `max_label_length = 239`.
- Эпоха обучения `epochs = 50`.
- Модель DenseNet.
- Размер входной свертки 280x32 и 3 компонентный.
- Оптимизатор Adam, скорость обучение 0.0005.
- Метрика для обучения CTC (Коннекционистская временная классификация) [20].

Модель DenseNet

Layer (type)	Output Shape	Param #	Connected to
input_data (InputLayer)	[(None, 32, 280, 3) 0]	[]	
zero_padding2d (ZeroPadding2D)	(None, 38, 286, 3) 0		['input_data[0][0]']
conv1/conv (Conv2D)	(None, 16, 140, 64) 9408		['zero_padding2d[0][0]']
conv1/bn (BatchNormalization)	(None, 16, 140, 64) 256		['conv1/conv[0][0]']
conv1/relu (Activation)	(None, 16, 140, 64) 0		['conv1/bn[0][0]']
zero_padding2d_1 (ZeroPadding2D)	(None, 18, 142, 64) 0		['conv1/relu[0][0]']
pool1 (MaxPooling2D)	(None, 8, 70, 64) 0		['zero_padding2d_1[0][0]']
conv2_block1_0_bn (BatchNormalization)	(None, 8, 70, 64) 256		['pool1[0][0]']
conv2_block1_0_relu (Activation)	(None, 8, 70, 64) 0		['conv2_block1_0_bn[0][0]']
conv2_block1_1_conv (Conv2D)	(None, 8, 70, 128) 8192		['conv2_block1_0_relu[0][0]']
conv2_block1_1_bn (BatchNormalization)	(None, 8, 70, 128) 512		['conv2_block1_1_conv[0][0]']
conv2_block1_1_relu (Activation)	(None, 8, 70, 128) 0		['conv2_block1_1_bn[0][0]']
conv2_block1_2_conv (Conv2D)	(None, 8, 70, 32) 36864		['conv2_block1_1_relu[0][0]']
conv2_block1_concat (Concatenate)	(None, 8, 70, 96) 0		['pool1[0][0]', 'conv2_block1_2_conv[0][0]']
conv2_block2_0_bn (BatchNormalization)	(None, 8, 70, 96) 384		['conv2_block1_concat[0][0]']
conv2_block2_0_relu (Activation)	(None, 8, 70, 96) 0		['conv2_block2_0_bn[0][0]']
conv2_block2_1_conv (Conv2D)	(None, 8, 70, 128) 12288		['conv2_block2_0_relu[0][0]']
conv2_block2_1_bn (BatchNormalization)	(None, 8, 70, 128) 512		['conv2_block2_1_conv[0][0]']
conv2_block2_1_relu (Activation)	(None, 8, 70, 128) 0		['conv2_block2_1_bn[0][0]']
conv2_block2_2_conv (Conv2D)	(None, 8, 70, 32) 36864		['conv2_block2_1_relu[0][0]']
conv2_block2_concat (Concatenate)	(None, 8, 70, 128) 0		['conv2_block1_concat[0][0]',

'conv2_block2_2_conv[0][0]'					
conv2_block3_0_bn (BatchNormalization)	(None,	8,	70,	128)	512
['conv2_block2_concat[0][0]'					
conv2_block3_0_relu (Activation)	(None, 8, 70, 128)	0			['conv2_block3_0_bn[0][0]'
conv2_block3_1_conv (Conv2D)	(None, 8, 70, 128)	16384			['conv2_block3_0_relu[0][0]'
conv2_block3_1_bn (BatchNormalization)	(None,	8,	70,	128)	512
['conv2_block3_1_conv[0][0]'					
conv2_block3_1_relu (Activation)	(None, 8, 70, 128)	0			['conv2_block3_1_bn[0][0]'
conv2_block3_2_conv (Conv2D)	(None, 8, 70, 32)	36864			['conv2_block3_1_relu[0][0]'
conv2_block3_concat (Concatenate)	(None, 8, 70, 160)	0			['conv2_block2_concat[0][0]'
'conv2_block3_2_conv[0][0]'					
conv2_block4_0_bn (BatchNormalization)	(None,	8,	70,	160)	640
['conv2_block3_concat[0][0]'					
conv2_block4_0_relu (Activation)	(None, 8, 70, 160)	0			['conv2_block4_0_bn[0][0]'
conv2_block4_1_conv (Conv2D)	(None, 8, 70, 128)	20480			['conv2_block4_0_relu[0][0]'
conv2_block4_1_bn (BatchNormalization)	(None,	8,	70,	128)	512
['conv2_block4_1_conv[0][0]'					
conv2_block4_1_relu (Activation)	(None, 8, 70, 128)	0			['conv2_block4_1_bn[0][0]'
conv2_block4_2_conv (Conv2D)	(None, 8, 70, 32)	36864			['conv2_block4_1_relu[0][0]'
conv2_block4_concat (Concatenate)	(None, 8, 70, 192)	0			['conv2_block3_concat[0][0]'
'conv2_block4_2_conv[0][0]'					
conv2_block5_0_bn (BatchNormalization)	(None,	8,	70,	192)	768
['conv2_block4_concat[0][0]'					
conv2_block5_0_relu (Activation)	(None, 8, 70, 192)	0			['conv2_block5_0_bn[0][0]'
conv2_block5_1_conv (Conv2D)	(None, 8, 70, 128)	24576			['conv2_block5_0_relu[0][0]'
conv2_block5_1_bn (BatchNormalization)	(None,	8,	70,	128)	512
['conv2_block5_1_conv[0][0]'					
conv2_block5_1_relu (Activation)	(None, 8, 70, 128)	0			['conv2_block5_1_bn[0][0]'
conv2_block5_2_conv (Conv2D)	(None, 8, 70, 32)	36864			['conv2_block5_1_relu[0][0]'
conv2_block5_concat (Concatenate)	(None, 8, 70, 224)	0			['conv2_block4_concat[0][0]'
'conv2_block5_2_conv[0][0]'					
conv2_block6_0_bn (BatchNormalization)	(None,	8,	70,	224)	896
['conv2_block5_concat[0][0]'					
conv2_block6_0_relu (Activation)	(None, 8, 70, 224)	0			['conv2_block6_0_bn[0][0]'
conv2_block6_1_conv (Conv2D)	(None, 8, 70, 128)	28672			['conv2_block6_0_relu[0][0]'
conv2_block6_1_bn (BatchNormalization)	(None,	8,	70,	128)	512
['conv2_block6_1_conv[0][0]'					
conv2_block6_1_relu (Activation)	(None, 8, 70, 128)	0			['conv2_block6_1_bn[0][0]'
conv2_block6_2_conv (Conv2D)	(None, 8, 70, 32)	36864			['conv2_block6_1_relu[0][0]'
conv2_block6_concat (Concatenate)	(None, 8, 70, 256)	0			['conv2_block5_concat[0][0]'
'conv2_block6_2_conv[0][0]'					
pool2_bn (BatchNormalization)	(None, 8, 70, 256)	1024			['conv2_block6_concat[0][0]'
pool2_relu (Activation)	(None, 8, 70, 256)	0			['pool2_bn[0][0]'
conv2d (Conv2D)	(None, 4, 70, 512)	3277312			['pool2_relu[0][0]'
reshape (Reshape)	(None, 280, 1, 512)	0			['conv2d[0][0]'
flatten (TimeDistributed)	(None, 280, 512)	0			['reshape[0][0]'
output (Dense)	(None, 280, 283)	145179			['flatten[0][0]'

==

Total params: 3,771,483

Trainable params: 3,767,579

Non-trainable params: 3,904

load model done.

```
sasha@sasha-ecoblue: /home/ram/2021-DIGIX-BASELINE/baseline-game5/recognizer

Total params: 3,771,483
Trainable params: 3,767,579
Non-trainable params: 3,904

Epoch 1/50
1/210 [.....] - ETA: 13:06 - loss: 1399.6299 - accuracy: 0.0000e+00
2/210 [.....] - ETA: 9:15 - loss: 795.0533 - accuracy: 0.0000e+00
3/210 [.....] - ETA: 9:12 - loss: 604.6071 - accuracy: 0.0000e+00
4/210 [.....] - ETA: 9:10 - loss: 519.9543 - accuracy: 0.0000e+00
5/210 [.....] - ETA: 9:08 - loss: 460.7217 - accuracy: 0.0000e+00
6/210 [.....] - ETA: 9:06 - loss: 418.5966 - accuracy: 0.0000e+00
7/210 [.....] - ETA: 9:03 - loss: 391.8034 - accuracy: 0.0000e+00
8/210 [.....] - ETA: 9:02 - loss: 369.0962 - accuracy: 0.0000e+00
9/210 [.....] - ETA: 8:59 - loss: 350.7556 - accuracy: 0.0000e+00
10/210 [.....] - ETA: 8:56 - loss: 335.1798 - accuracy: 0.0000e+00
11/210 [.....] - ETA: 8:53 - loss: 320.2180 - accuracy: 0.0000e+00
12/210 [.....] - ETA: 8:50 - loss: 310.4994 - accuracy: 0.0000e+00
13/210 [.....] - ETA: 8:48 - loss: 300.6336 - accuracy: 0.0000e+00
14/210 [.....] - ETA: 8:46 - loss: 293.6300 - accuracy: 0.0000e+00
15/210 [.....] - ETA: 8:43 - loss: 287.5651 - accuracy: 0.0000e+00
16/210 [.....] - ETA: 8:39 - loss: 280.3647 - accuracy: 0.0000e+00
17/210 [.....] - ETA: 8:37 - loss: 273.5972 - accuracy: 0.0000e+00
18/210 [.....] - ETA: 8:35 - loss: 268.0829 - accuracy: 0.0000e+00
19/210 [.....] - ETA: 8:32 - loss: 262.6168 - accuracy: 0.0000e+00
20/210 [.....] - ETA: 8:29 - loss: 257.7321 - accuracy: 0.0000e+00
21/210 [.....] - ETA: 8:27 - loss: 252.7703 - accuracy: 0.0000e+00
22/210 [.....] - ETA: 8:24 - loss: 247.5339 - accuracy: 0.0000e+00
23/210 [.....] - ETA: 8:21 - loss: 243.1859 - accuracy: 0.0000e+00
24/210 [.....] - ETA: 8:18 - loss: 238.1078 - accuracy: 0.0000e+00
25/210 [.....] - ETA: 8:15 - loss: 232.2779 - accuracy: 0.0000e+00
26/210 [.....] - ETA: 8:13 - loss: 228.0491 - accuracy: 0.0000e+00
27/210 [.....] - ETA: 8:10 - loss: 224.4259 - accuracy: 0.0000e+00
28/210 [.....] - ETA: 8:07 - loss: 221.4850 - accuracy: 0.0000e+00
29/210 [.....] - ETA: 8:04 - loss: 218.2506 - accuracy: 0.0000e+00
30/210 [.....] - ETA: 8:02 - loss: 214.7959 - accuracy: 0.0000e+00
31/210 [.....] - ETA: 7:59 - loss: 211.4659 - accuracy: 0.0000e+00
32/210 [.....] - ETA: 7:56 - loss: 208.2757 - accuracy: 0.0000e+00
33/210 [.....] - ETA: 7:54 - loss: 206.5795 - accuracy: 0.0000e+00
34/210 [.....] - ETA: 7:51 - loss: 204.0060 - accuracy: 0.0000e+00
35/210 [.....] - ETA: 7:49 - loss: 201.0288 - accuracy: 0.0000e+00
36/210 [.....] - ETA: 7:46 - loss: 198.8258 - accuracy: 0.0000e+00
37/210 [.....] - ETA: 7:44 - loss: 196.6810 - accuracy: 0.0000e+00
38/210 [.....] - ETA: 7:41 - loss: 194.8691 - accuracy: 0.0000e+00
39/210 [.....] - ETA: 7:38 - loss: 192.6026 - accuracy: 0.0000e+00
40/210 [.....] - ETA: 7:36 - loss: 189.8880 - accuracy: 0.0000e+00
41/210 [.....] - ETA: 7:33 - loss: 187.5860 - accuracy: 0.0000e+00
42/210 [.....] - ETA: 7:30 - loss: 185.7580 - accuracy: 0.0000e+00
43/210 [.....] - ETA: 7:28 - loss: 183.6941 - accuracy: 0.0000e+00
44/210 [.....] - ETA: 7:25 - loss: 182.0972 - accuracy: 0.0000e+00
45/210 [.....] - ETA: 7:22 - loss: 180.7160 - accuracy: 0.0000e+00
46/210 [.....] - ETA: 7:20 - loss: 179.1184 - accuracy: 0.0000e+00
47/210 [.....] - ETA: 7:17 - loss: 177.5343 - accuracy: 0.0000e+00
48/210 [.....] - ETA: 7:14 - loss: 175.8707 - accuracy: 0.0000e+00
49/210 [.....] - ETA: 7:12 - loss: 174.2960 - accuracy: 0.0000e+00
50/210 [.....] - ETA: 7:09 - loss: 172.6408 - accuracy: 0.0000e+00
51/210 [.....] - ETA: 7:06 - loss: 171.0877 - accuracy: 0.0000e+00
```

Рисунок 17. Выполнение программы обучение модели

Процесс аугментация данными:

1. Изменение размер изображение на 280x32
2. Преобразования в тензор и нормирования от -1 до 1

Таблица 5. Вывод процесс обучение модели

Epoch (Эпоха)	Loss (CTC)	accuracy	Lr (Скорость)	Время обучения (секунда)
1	102.3565	0.000000	5.000000e-04	569
2	70.6220	0.000000	5.000000e-04	569
3	54.2058	0.000075	5.000000e-04	573
4	40.1689	0.002700	5.000000e-04	573
5	34.6726	0.018500	5.000000e-04	571
6	31.1881	0.052000	5.000000e-04	570

Epoch (Эпоха)	Loss (CTC)	accuracy	Lr (Скорость)	Время обучения (секунда)
7	29.1256	0.081900	5.000000e-04	570
8	27.6356	0.106600	5.000000e-04	570
9	26.2397	0.128200	5.000000e-04	574
10	25.3960	0.146100	5.000000e-04	581
11	24.2791	0.162600	5.000000e-04	569
12	23.3033	0.178700	5.000000e-04	567
13	22.6908	0.195900	5.000000e-04	575
14	22.1181	0.201400	5.000000e-04	574
15	21.4371	0.216200	5.000000e-04	571
16	21.0644	0.226400	5.000000e-04	571
17	20.2271	0.237100	5.000000e-04	565
18	20.0330	0.242500	5.000000e-04	562
19	19.6939	0.255000	5.000000e-04	562
20	18.6772	0.264600	5.000000e-04	562
21	18.3791	0.270700	5.000000e-04	562
22	18.0687	0.279600	5.000000e-04	562
23	17.8101	0.284900	5.000000e-04	563
24	17.0091	0.296600	5.000000e-04	580
25	16.7885	0.302800	5.000000e-04	578
26	16.5632	0.310300	5.000000e-04	583
27	15.6259	0.322900	5.000000e-04	589
28	15.7555	0.329400	5.000000e-04	594
29	13.1467	0.413100	5.000000e-05	593
30	12.2423	0.441500	5.000000e-05	592
31	12.3565	0.444700	5.000000e-05	597
32	12.1558	0.455000	5.000000e-06	605
33	11.7739	0.462700	5.000000e-06	589
34	12.0154	0.459100	5.000000e-06	585
35	11.7770	0.462300	5.000000e-07	594
36	11.8297	0.461500	5.000000e-08	580
37	11.9964	0.463500	5.000000e-09	582

Epoch (Эпоха)	Loss (CTC)	accuracy	Lr (Скорость)	Время обучения (секунда)
38	11.6463	0.461100	5.000000e-10	597
39	12.0612	0.461600	5.000000e-10	583
40	11.9004	0.464100	5.000000e-11	590
41	11.8442	0.459000	5.000000e-12	589
42	11.8288	0.465000	5.000000e-13	584
43	11.5611	0.461900	5.000000e-14	580
44	12.0484	0.462600	5.000000e-14	578
45	11.8042	0.464700	5.000000e-15	572
46	11.9628	0.459400	5.000000e-16	570
47	11.6019	0.461400	5.000000e-17	569
48	11.8933	0.463800	5.000000e-18	570

Всего время обучения: 28 335 секунда (7 часов, 52 минута, 15 секунда).

Время выполнения обучения 27 708 секунда (7 часов, 41 минута, 48 секунда).

Средняя время обучения за одну эпоху: 577.25 (9 минут, 37.25 секунда).

Средняя время шаг обучения 2,699 секунда (3 секунда).

Использовалось CPU без графического ускорения, т.к. для графического ускорения не хватало памяти для выполнения обучения.

На этапе обучение эпоха 48 завершилось, т.к. функция потери перестала уменьшаться, произошло затухание градиента. Сработал критерий останова (табл. 5).

```
sasha@sasha-ecoblue: /home/ram/2021-DIGIX-BASELINE/baseline-game5/recognizer
156/210 [====>.....] - ETA: 2:25 - loss: 12.3053 - accuracy: 0.4631
157/210 [====>.....] - ETA: 2:22 - loss: 12.3200 - accuracy: 0.4630
158/210 [====>.....] - ETA: 2:20 - loss: 12.2850 - accuracy: 0.4631
159/210 [====>.....] - ETA: 2:17 - loss: 12.2869 - accuracy: 0.4638
160/210 [====>.....] - ETA: 2:14 - loss: 12.2823 - accuracy: 0.4637
161/210 [====>.....] - ETA: 2:12 - loss: 12.3082 - accuracy: 0.4641
162/210 [====>.....] - ETA: 2:09 - loss: 12.2797 - accuracy: 0.4647
163/210 [====>.....] - ETA: 2:06 - loss: 12.2695 - accuracy: 0.4647
164/210 [====>.....] - ETA: 2:03 - loss: 12.2530 - accuracy: 0.4646
165/210 [====>.....] - ETA: 2:01 - loss: 12.2271 - accuracy: 0.4643
166/210 [====>.....] - ETA: 1:58 - loss: 12.2101 - accuracy: 0.4640
167/210 [====>.....] - ETA: 1:55 - loss: 12.1853 - accuracy: 0.4640
168/210 [====>.....] - ETA: 1:53 - loss: 12.2475 - accuracy: 0.4634
169/210 [====>.....] - ETA: 1:50 - loss: 12.2153 - accuracy: 0.4637
170/210 [====>.....] - ETA: 1:47 - loss: 12.1941 - accuracy: 0.4635
171/210 [====>.....] - ETA: 1:45 - loss: 12.1702 - accuracy: 0.4635
172/210 [====>.....] - ETA: 1:42 - loss: 12.1615 - accuracy: 0.4634
173/210 [====>.....] - ETA: 1:39 - loss: 12.1183 - accuracy: 0.4638
174/210 [====>.....] - ETA: 1:37 - loss: 12.0971 - accuracy: 0.4638
175/210 [====>.....] - ETA: 1:34 - loss: 12.0799 - accuracy: 0.4638
176/210 [====>.....] - ETA: 1:31 - loss: 12.1158 - accuracy: 0.4636
177/210 [====>.....] - ETA: 1:28 - loss: 12.1076 - accuracy: 0.4639
178/210 [====>.....] - ETA: 1:26 - loss: 12.0966 - accuracy: 0.4636
179/210 [====>.....] - ETA: 1:23 - loss: 12.0786 - accuracy: 0.4636
180/210 [====>.....] - ETA: 1:20 - loss: 12.0418 - accuracy: 0.4635
181/210 [====>.....] - ETA: 1:18 - loss: 12.0297 - accuracy: 0.4631
182/210 [====>.....] - ETA: 1:15 - loss: 12.0163 - accuracy: 0.4633
183/210 [====>.....] - ETA: 1:12 - loss: 12.0181 - accuracy: 0.4634
184/210 [====>.....] - ETA: 1:10 - loss: 12.0105 - accuracy: 0.4634
185/210 [====>.....] - ETA: 1:07 - loss: 12.0035 - accuracy: 0.4630
186/210 [====>.....] - ETA: 1:04 - loss: 12.0271 - accuracy: 0.4630
187/210 [====>.....] - ETA: 1:02 - loss: 12.0242 - accuracy: 0.4632
188/210 [====>.....] - ETA: 59s - loss: 12.0256 - accuracy: 0.4627
189/210 [====>.....] - ETA: 56s - loss: 12.0193 - accuracy: 0.4628
190/210 [====>.....] - ETA: 53s - loss: 12.0085 - accuracy: 0.4626
191/210 [====>.....] - ETA: 51s - loss: 11.9813 - accuracy: 0.4628
192/210 [====>.....] - ETA: 48s - loss: 11.9726 - accuracy: 0.4629
193/210 [====>.....] - ETA: 45s - loss: 11.9798 - accuracy: 0.4629
194/210 [====>.....] - ETA: 43s - loss: 11.9916 - accuracy: 0.4631
195/210 [====>.....] - ETA: 40s - loss: 11.9930 - accuracy: 0.4630
196/210 [====>.....] - ETA: 37s - loss: 12.0009 - accuracy: 0.4629
197/210 [====>.....] - ETA: 35s - loss: 11.9788 - accuracy: 0.4629
198/210 [====>.....] - ETA: 32s - loss: 11.9634 - accuracy: 0.4630
199/210 [====>.....] - ETA: 29s - loss: 11.9509 - accuracy: 0.4630
200/210 [====>.....] - ETA: 26s - loss: 11.9538 - accuracy: 0.4632
201/210 [====>.....] - ETA: 24s - loss: 11.9325 - accuracy: 0.4632
202/210 [====>.....] - ETA: 21s - loss: 11.9210 - accuracy: 0.4635
203/210 [====>.....] - ETA: 18s - loss: 11.8971 - accuracy: 0.4638
204/210 [====>.....] - ETA: 16s - loss: 11.9113 - accuracy: 0.4635
205/210 [====>.....] - ETA: 13s - loss: 11.8974 - accuracy: 0.4635
206/210 [====>.....] - ETA: 10s - loss: 11.9138 - accuracy: 0.4639
207/210 [====>.....] - ETA: 8s - loss: 11.8968 - accuracy: 0.4636
208/210 [====>.....] - ETA: 5s - loss: 11.9024 - accuracy: 0.4638
209/210 [====>.....] - ETA: 2s - loss: 11.9006 - accuracy: 0.4639
210/210 [====>.....] - ETA: 0s - loss: 11.8933 - accuracy: 0.4638
Epoch 48: early stopping
user@user-studio-3d: /home/ram/2021-DIGIX-BASELINE/baseline-game5/recognizer$
```

Рисунок 18. Программа завершила обучение, так как скорость обучения упала, градиент затух

Loss — функция потерь, метрика СТС.

Accuracy — оценка обученности модели на тренировочных данных.

Lr — скорость обучения с использованием оптимизатора Adam.

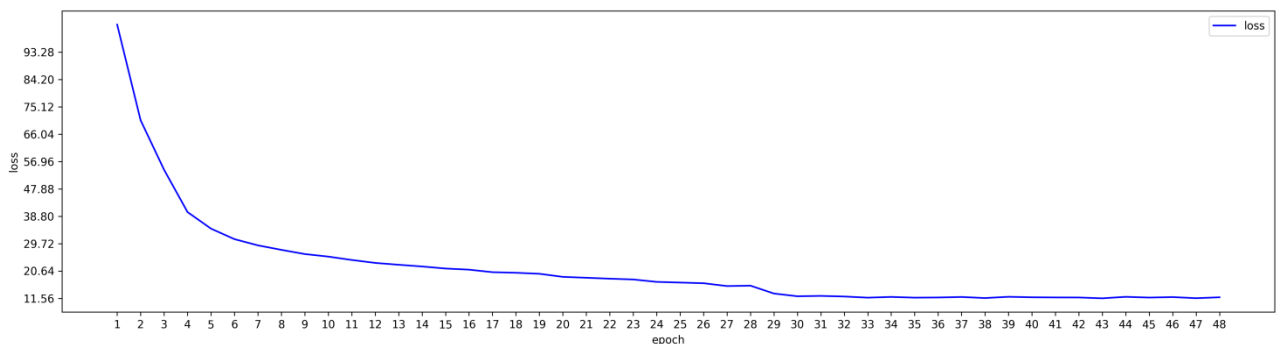


Рисунок 19. На графике представлены по оси x — epoch, по оси y — loss

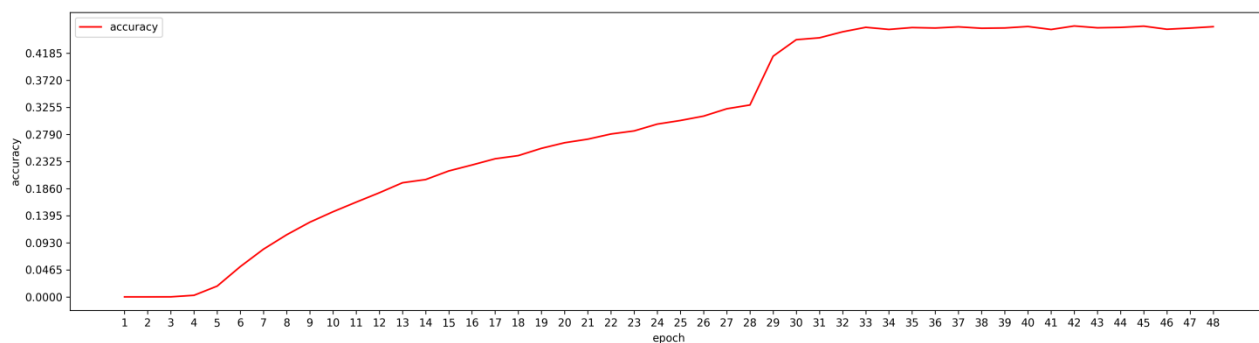


Рисунок 20. На графике представлены по оси x — epoch, по оси y — accuracy

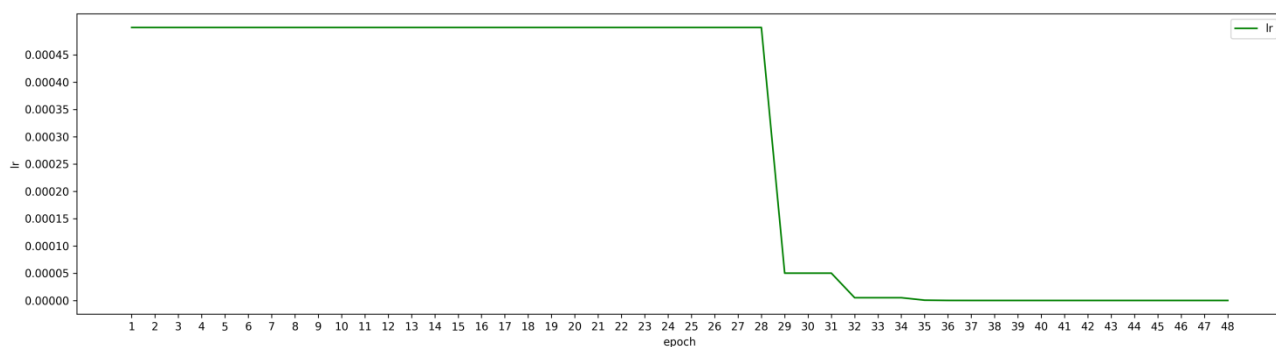


Рисунок 21. На графике представлены по оси x — epoch, по оси y — lr

На графике видно, что (Рис. 19) начиная с эпохи 28, качество модели резко улучшилось, а дальше держится примерно 0.46, а на другом графике (Рис 18), функция потерь держится примерно 11. На графике (Рис 20) скорость обучение начала падать, начиная с эпохи 29 и постепенно стало быстрее падать, т.к. модель стала хуже обучается.

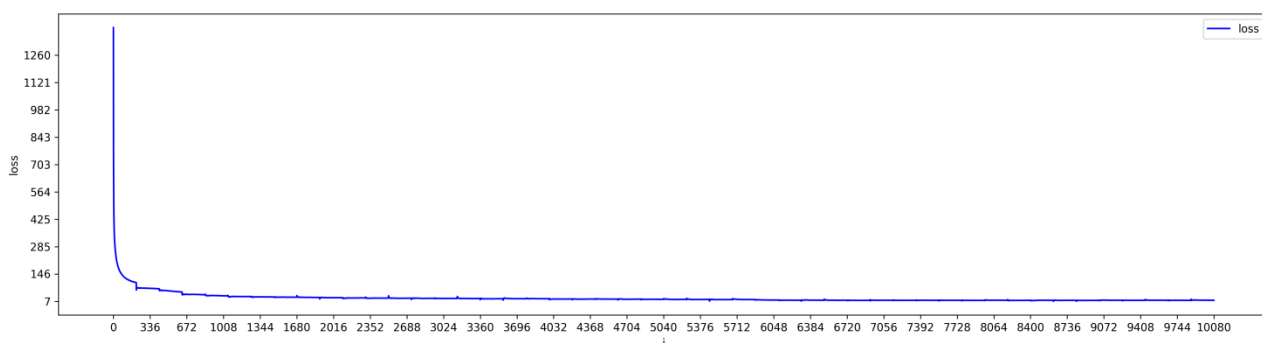


Рисунок 22. На графике представлены по оси x — цикл обучения, по оси y — loss. За весь период обучения



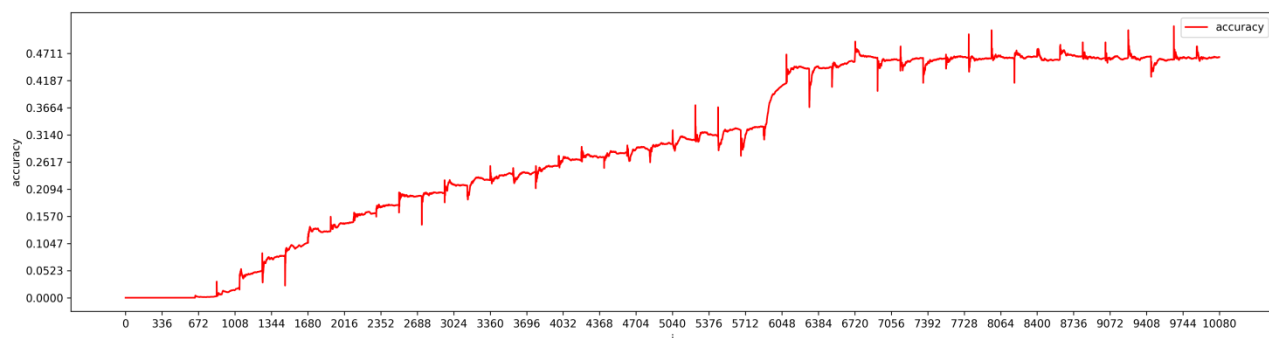


Рисунок 23. На графике представлены по оси x — цикл обучения, по оси y — accuracy. За весь период обучения

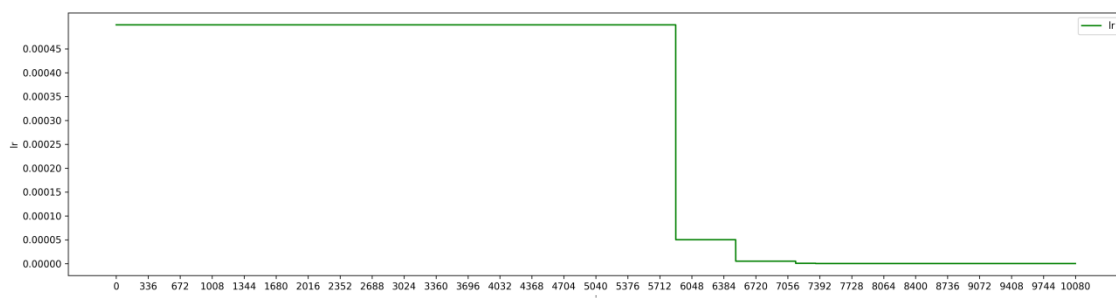


Рисунок 24. На графике представлены по оси x — цикл обучения, по оси y — lr. За весь период обучения

### 3.4. Запуск detector/predict.py

В программе выполняется предсказание модели по поиску текста в изображении (Detector).

Этапы подготовка изображения:

3. Изменяем размер изображения в приближенный 1920 пикселей.
4. Преобразуем в Tensor и RGB каналы.
5. Нормализация.



Рисунок 25. Сравнение качество обучение, слева — исходные разметки, а справа — разметка результат - обученный модель

Средняя время выполнение обработки одно изображение (учитывая распознавание, ввод модели, рисование область в изображение, вывод в файл) равно 3.601 секунде.

Время выполнение программы для всего изображения 14958.781 секунда (4 часа, 9 минута, 18.781 секунда).

### 3.5. Запуск recognizer/predict.py

В программе выполняется предсказание модели по извлечению текста в изображения (Recognizer).

```

user@user-studio-3D: /mnt/Data/Project/Python/AI/Milfslm/App/recognizer
Rd
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 18ms/step
e
1/1 [=====] - 0s 24ms/step
1/1 [=====] - 0s 19ms/step
2.70G
1/1 [=====] - 0s 19ms/step
1/1 [=====] - 0s 19ms/step
TAI
1/1 [=====] - 0s 18ms/step
YA
1/1 [=====] - 0s 19ms/step
TA
1/1 [=====] - 0s 18ms/step
Dna
1/1 [=====] - 0s 19ms/step
onol
1/1 [=====] - 0s 18ms/step

```

Рисунок 26. Выполнение программы

```
IPython: App/detector
{'label': 'CHDPTU5L',
 'points': [['363', '314'], ['545', '312'], ['545', '369'], ['363', '371']]},
{'label': '',
 'points': [['592', '301'], ['599', '296'], ['603', '302'], ['596', '307']]},
{'label': 'AU',
 'points': [['549', '274'], ['598', '268'], ['603', '310'], ['554', '316']]},
{'label': 'CENTRE',
 'points': [['594', '263'], ['746', '263'], ['746', '320'], ['594', '320']]},
{'label': 'APERTU DE LA CASA',
 'points': [['502', '175'], ['791', '174'], ['791', '222'], ['502', '223']]},
{'label': '',
 'points': [['894', '92'], ['915', '91'], ['916', '105'], ['895', '106']]},
{'label': '',
 'points': [['510', '98'], ['513', '79'], ['524', '81'], ['521', '100']]},
{'label': 'la',
 'points': [['895', '63'], ['918', '63'], ['918', '80'], ['895', '80']]},
{'label': '',
 'points': [['784', '108'], ['781', '63'], ['817', '60'], ['820', '105']]},
{'label': 'DTA',
 'points': [['803', '56'], ['900', '56'], ['900', '113'], ['803', '113']]},
{'label': 't',
 'points': [['500', '63'], ['519', '65'], ['517', '80'], ['498', '77']]},
{'label': 'ORIGENS',
 'points': [['521', '51'], ['786', '51'], ['786', '117'], ['521', '117']]},
```

Рисунок 2.27. Пример чтение файла test\_0.json

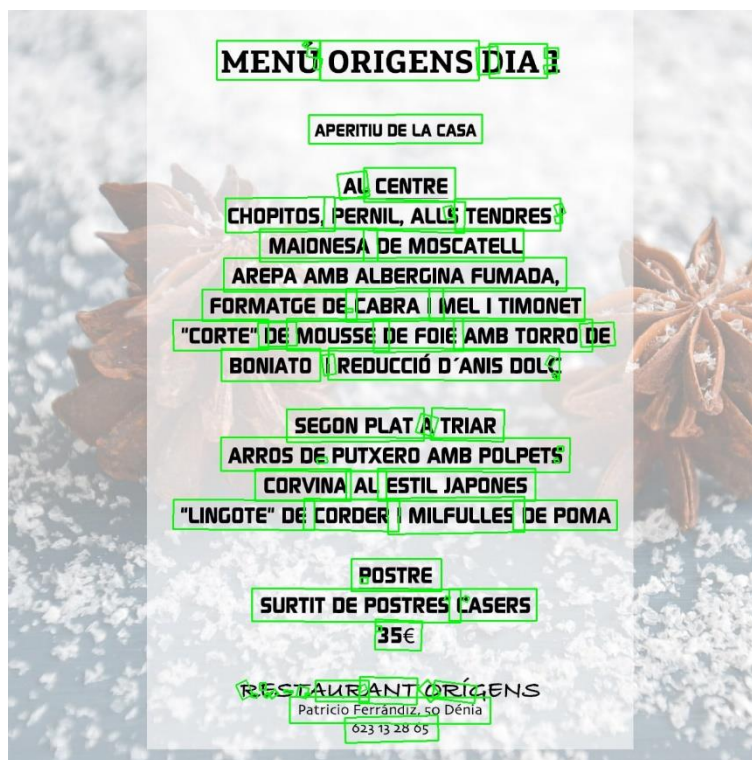


Рисунок 2.28. Файл test\_0.JPEG. Выделенное квадратиком текста

Время выполнение программы равно 4292,502 секунд (1 час, 11 минут, 32.502 секунда).

Средняя время выполнение обработки каждого изображения составило 1.254 секунд.

Всего обработано изображений из 3423 файлов.

## Заклучение

В данной работе представлена разработанная программа для поиска текста в изображениях и получения текста из изображений. В результате обучена модель для оптического распознавания текста. В ходе обучения определена метрика для оценки модели. Несмотря на ошибки и неточности в ходе обучения модели искусственного интеллекта получилось неплохо, модель может при определенных входных данных выдать правильные предсказания. В данном случае модель машинного обучения показала хороший результат по поиску текста изображения с точностью 83,52%. Слабее получен результат модели по извлечению текста из изображения - результат с точностью 46,38%. Время выполнения программы составляет 79171,538 секунд (21 час, 59 минут, 31.538 секунда) (табл 6).

Таблица 6. Время выполнения программы

Программа	Время
detector/train.py	32136.69 секунд (8 часов, 55.61 минут, 36.69 секунда)
Подготовка данных к обучению recognizer	1 минута и 15,565 секунда.
recognizer/train.py	27708 секунд (7 часов, 41 минута, 48 секунда).
detector/predict.py	14958.781 секунд (4 часа, 9 минута, 18.781 секунда).
recognizer/predict.py	4292,502 секунд (1 час, 11 минут, 32.502 секунда).
Сумма	79171,538 секунд (21 час, 59 минут, 31.538 секунда)

## Список использованной литературы

1. Culjak I. et al. A brief introduction to OpenCV //2012 proceedings of the 35th international convention MIPRO. – IEEE, 2012. – С. 1725-1730.
2. Pulli K. et al. Real-time computer vision with OpenCV //Communications of the ACM. – 2012. – Т. 55. – №. 6. – С. 61-69.
3. Xie G., Lu W. Image edge detection based on opencv //International Journal of Electronics and Electrical Engineering. – 2013. – Т. 1. – №. 2. – С. 104-106.
4. Agam G. Introduction to programming with OpenCV //Online Document. – 2006. – Т. 27.
5. Qiao L. et al. DavarOCR: A Toolbox for OCR and Multi-Modal Document Understanding //Proceedings of the 30th ACM International Conference on Multimedia. – 2022. – С. 7355-7358.
6. Visintin M., Dever A. Optical Character Recognition. Master Degree Thesis. 2018.
7. Lu N. et al. Master: Multi-aspect non-local network for scene text recognition //Pattern Recognition. – 2021. – Т. 117. – С. 107980.
8. Yu W. et al. PICK: processing key information extraction from documents using improved graph learning-convolutional networks //2020 25th International Conference on Pattern Recognition (ICPR). – IEEE, 2021. – С. 4363-4370.
9. Lin P. H., Wang D. A Optical Character Recognition System with Pattern Editing Functionalities.
10. Kuang Z. et al. MMOCR: a comprehensive toolbox for text detection, recognition and understanding //Proceedings of the 29th ACM International Conference on Multimedia. – 2021. – С. 3791-3794.
11. Zayene O. et al. ICPR2020 competition on text detection and recognition in arabic news video frames //Pattern Recognition. ICPR International Workshops and Challenges: Virtual Event, January 10-15, 2021, Proceedings, Part VIII. – Springer International Publishing, 2021. – С. 343-356.

12. Jaume G., Ekenel H. K., Thiran J. P. Funsd: A dataset for form understanding in noisy scanned documents //2019 International Conference on Document Analysis and Recognition Workshops (ICDARW). – IEEE, 2019. – Т. 2. – С. 1-6.
13. Batra P. et al. A novel memory and time-efficient ALPR system based on YOLOv5 //Sensors. – 2022. – Т. 22. – №. 14. – С. 5283.
14. Pamnani M. DOT-HAZMAT: An Android application for detection of hazardous materials (HAZMAT) A student article showcasing an Android application called DOT-HAZMAT based on Machine Learning and Computer Vision that locates and categorizes HAZMAT placards in perilous rescue operations.
15. Vesalainen A. Image Segmentation methods for fine-grained OCR Document Layout Analysis. Master's thesis. 2022.
16. Olejniczak K., Šulc M. Text Detection Forgot About Document OCR //arXiv preprint arXiv:2210.07903. 2022.
17. GLOBAL AI CHALLENGE URL:  
<https://developer.huawei.com/consumer/en/activity/devStarAI/algo/> (дата обращения: 2021-10-04).
18. GitHub - luoda888/2021-DIGIX-BASELINE: 2021 huawei DIGIX competition baseline // GitHub URL: <https://github.com/luoda888/2021-DIGIX-BASELINE> (дата обращения: 2022-05-26).
19. Dice — PyTorch-Metrics 0.11.4 documentation // Dice URL:  
<https://torchmetrics.readthedocs.io/en/stable/classification/dice.html> (дата обращения: 2023-06-27).
20. CTCLoss — PyTorch 2.0 documentation // CTCLoss URL:  
<https://pytorch.org/docs/stable/generated/torch.nn.CTCLoss.html> (дата обращения: 2023-06-27).