







## **Анотація**

Робота присвячена розробці системи розпізнавання тексту для автоматизації обробки документів. У століття цифрової трансформації все більше сфер бізнесу починають використовувати сучасні технології розпізнавання тексту для автоматизації бізнес процесів і підвищення якості обслуговування клієнтів. У даній роботі досліджуються існуючі технології оптичного розпізнавання тексту, такі як: Google Vision, Tesseract і інші. Досліджено методи попередньої обробки зображень і пост-обробки результатів для поліпшення якості розпізнавання тексту. Розроблена система розпізнавання тексту з фотографій закордонних паспортів. Розроблена система дозволяє автоматизувати обробку документів і в перспективі здатна заощадити безліч людських ресурсів.

## **Abstract**

The work is devoted to the development of a text recognition system for automating the processing of documents. In the age of digital transformation, more and more business areas are starting to use modern text recognition technologies to automate business processes and improve the quality of customer service. This paper explores existing optical text recognition technologies, such as Google Vision, Tesseract, and others. The methods of image preprocessing and post-processing of the results are studied to improve the quality of text recognition. A system for recognizing text from photographs of passports has been implemented. The implemented system allows you to automate the processing of documents and in the future can save many human resources.







- Мати змогу оброблювати фотографії, зроблені на непрофесійну фотокамеру;
- Показувати високу точність розпізнавання тексту.

### **5.2. Вимоги до програмного забезпечення**

- Операційна система MS Windows XP, MS Windows Vista, MS Windows 7, MS Windows 8/8.1, MS Windows 10;
- Python 3 і вище;
- Telegram 2.1. і вище

### **5.3. Вимоги до апаратного забезпечення**

- Комп'ютер на базі процесору Intel Pentium 2 і вище
- Оперативної пам'яті не менше 128 Мбайт

					<i>ІАЛЦ.467100.002 ТЗ</i>	Арк.
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>		3

# Пояснювальна записка до дипломного проекту

на тему: Система розпізнавання тексту на документах

Київ - 2020 року





5.1.1.	Оцінка роботи системи для паспортів різних країн .....	47
5.1.2.	Оцінка роботи системи для фотографій, зроблених під кутом .....	49
4.1.3.	Оцінка роботи системи для фотографій, зроблених при поганому освітленні .....	51
5.2.	Аналіз результатів експериментів .....	53
ВИСНОВКИ ДО РОЗДІЛУ 4 .....		55
ВИСНОВКИ.....		56
ПЕРЕЛІК ПОСИЛАНЬ.....		57
ДОДАТКИ.....		60



### **Практичне значення**

Розроблена система, використовуючи передові технології розпізнавання тексту, здатна автоматизувати рутинну роботу по обробці даних з закордонних паспортів. Ця система може бути використана у туристичних компаніях чи державних установах для автоматизації обслуговування клієнтів та економії часових та людських ресурсів.

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		6















маркування зображення, розпізнавання обличчя та оптичне розпізнавання символів (OCR) [11].

Розпізнавання тексту виконує оптичне розпізнавання символів (OCR). Технологія виявляє та витягує текст із зображення із підтримкою широкого кола мов. Vision API також має автоматичну ідентифікацію мови.

Послуга розпізнавання тексту за допомогою Google Cloud Vision API є платною, але в місяць платформа дозволяє проводити 1000 API безкоштовних запитів на розпізнавання тексту, далі – 1,5\$ за кожні 1000 запитів [11].

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		14

## ВИСНОВКИ ДО РОЗДІЛУ 1

1. Аналіз сучасних технологій оптичного розпізнавання тексту показав, що для обраної задачі розпізнавання паспортів найбільш доцільно використовувати наступні технології:
  - Google Cloud Vision API
  - Tesseract OCR
2. Обидві обрані технології мають свої переваги та недоліки.

Google Cloud Vision API має зручний у використанні інтерфейс та показує високу якість розпізнавання тексту, але є платним (\$1.5 за 1000 запитів) та для роботи потребує доступ до інтернету.

Tesseract OCR є безкоштовним фреймворком з відкритим кодом, для якого вже існують навчання моделі для багатьох мов. Tesseract OCR не потребує доступу до інтернету та здатний працювати навіть на мобільних пристроях. Також, tesseract дозволяє перенавчати існуючі моделі для своїх даних та робити детальне налаштування для покращення роботи моделі. Недоліком можна зазначити більш складний інтерфейс для взаємодії та меншу точність розпізнавання тексту в порівнянні з Google Cloud Vision API.

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		15























[https://api.telegram.org/file/bot<telegram\\_bot\\_token>/<path\\_to\\_file>](https://api.telegram.org/file/bot<telegram_bot_token>/<path_to_file>), де з відповіді взято < path\_to\_file>. Гарантується, що посилання буде дійсною не менше 1 години. Коли посилання закінчується, нове можна отримати, надіславши запит *getFile* ще раз.

## **2.4. Python-telegram-bot**

Для більш зручної та ефективної розробки телеграм ботів існують багато бібліотек для роботи з Telegram bot API. У цій роботі ми будемо використовувати найвідомішу з них - python-telegram-bot [21].

Бібліотека python-telegram-bot забезпечує чистий інтерфейс Python для API Telegram Bot . Він сумісний з версіями Python 3.5+ та PyPy . Всі типи та методи API Telegram Bot API 4.8 підтримуються [22].

Окрім чистої реалізації API, у цій бібліотеці є ряд класів високого рівня, щоб зробити розробку ботів легкою та зрозумілою. Ці класи містяться в telegram.ext модулі [22].

					ІАЛЦ.467100.003 ПЗ	Арк.
						27
Зм.	Арк.	№ докум.	Підпис	Дата		

## ВИСНОВКИ ДО РОЗДІЛУ 2

1. Аналіз засобів для розробки системи оптичного розпізнавання тексту на фотографіях закордонних паспортів показав, що найбільш доцільно використовувати Google Vision API у якості технології для розпізнавання тексту та Telegram Bot API у якості зручного графічного інтерфейсу користувача;
2. Обрані технології гарантують надійну та безперебійну роботу програми та безпеку персональних даних;
3. В якості мови програмування була обрана мова Python, оскільки вона має підтримку бібліотек для роботи з Google Cloud Platform та Telegram Bot API.

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		28































Розгорнувши Flask API для розпізнавання MRZ Heroку, та запустивши додаток під назвою *pass-bot* ми можемо надсилати запити до API на наступний url: <https://pass-bot.herokuapp.com/api>.

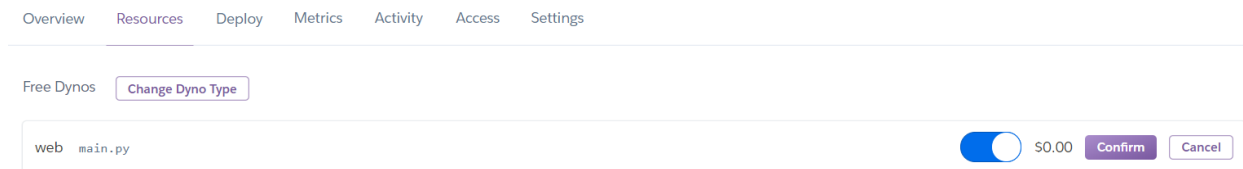


Рис.3.13. Запуск додатку на Heroku

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		44

### ВИСНОВКИ ДО РОЗДІЛУ 3

1. В даному розділі було розроблено систему розпізнавання тексту на документах, а саме розпізнавання MRZ коду з закордонних паспортів;
2. Був розроблений інтерфейс користувача для зручної взаємодії з системою;
3. Був розроблений алгоритм розпізнавання тексту з зображення та алгоритм автокорекції помилок;
4. Було розроблено сценарій взаємодії користувача із програмою за допомогою телеграм боту;
5. Було розроблено структурну взаємодії графічного інтерфейсу з серверною частиною програми.

					<i>ІАЛЦ.467100.003 ПЗ</i>	Арк.
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>		45

















- Система здатна оброблювати не лише якісні зображення, зроблені за допомогою сканера, а й фотознімки низької якості, зроблені за допомогою фотокамери користувача;
- Система здатна розпізнавати текст з зображень, зроблених під кутом та при поганому освітленні, однак в деяких випадках це може вплинути на якість розпізнавання. Тому, слід явно вказувати користувачу у рекомендації як правильно зробити фото паспорту, задля досягнення максимальної точності розпізнавання.

Підсумовуючи усі перераховані висновки, можна сказати, що розроблена система розпізнавання тексту вже на даному етапі здатна автоматизувати та значно скоротити час на обробку фотознімків закордонних паспортів, а це в свою чергу надає велику цінність продукту для бізнесу.

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		54

## ВИСНОВКИ ДО РОЗДІЛУ 4

1. В даному розділі було представлено тестування роботи системи оптичного розпізнавання тексту з закордонних паспортів;
2. Були проведені експерименти для оцінки якості розпізнавання тексту, а саме: тестування на паспортах різних країн, тестування на фотознімках паспортів під різним кутом нахилу та при поганому освітленні;
3. Експериментально було виявлено чинники, які впливають на якість розпізнавання тексту, а саме: чіткість зображення, кут нахилу тексту, освітлення та наявність бліків на зображенні;
4. Також, було доведено, що система здатна працювати не тільки із якісними зображеннями, а й з реальними фотознімками, зробленими користувачем на камеру телефону при поганому освітленні.

					<i>ІАЛЦ.467100.003 ПЗ</i>	Арк.
						55
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>		

## ВИСНОВКИ

1. У представленій роботі було розглянуто методи оптичного розпізнавання тексту. Були розглянуті передові технології в області розпізнавання тексту. Було аргументовано вибір технології Google Vision API для задачі розпізнавання тексту на закордонних паспортах.
2. Буда розроблена система оптичного розпізнавання тексту з закордонних паспортів та графічний інтерфейс до неї у вигляді Telegram Bot.
3. Було проведено тестування роботи системи та аналіз результатів за допомогою ряду експериментів. Оцінка якості роботи показала, що система здатна опрацьовувати фотознімки паспортів різної якості, у деяких випадках, навіть при поганому освітленні та наявному тексті під кутом нахилу.

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		56

## ПЕРЕЛІК ПОСИЛАНЬ

1. Что такое OCR [Електронний ресурс]. – АБВУ FineReader, 2019. – Режим доступу: <https://www.abby.com/ru-ru/finereader/what-is-ocr/>.
2. Закордонний паспорт громадянина України для виїзду за кордон [Електронний ресурс]. – ДП "Документ", 2020, – Режим доступу: <https://rivne.pasport.org.ua/poslugi/zakordonnij-pasport>.
3. Біометричні документи в Україні [Електронний ресурс]. – ДМС України, 2020, – Режим доступу: <https://dmsu.gov.ua/faq/biometriczni-dokumenti-v-ukrajni.html>.
4. Паспорт громадянина України для виїзду за кордон [Електронний ресурс]. – Cycolwiki, 2020, – Режим доступу: [http://cyclowiki.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BE%D1%87%D0%B8%D1%82%D0%B0%D0%B5%D0%BC%D1%8B%D0%B9\\_%D0%BF%D0%B0%D1%81%D0%BF%D0%BE%D1%80%D1%82](http://cyclowiki.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BE%D1%87%D0%B8%D1%82%D0%B0%D0%B5%D0%BC%D1%8B%D0%B9_%D0%BF%D0%B0%D1%81%D0%BF%D0%BE%D1%80%D1%82).
5. АБВУ. What is OCR and OCR technology [Електронний ресурс] / АБВУ – Режим доступу до ресурсу: <https://www.abby.com/en-eu/finereader/what-is-ocr/>.
6. OCR System FineReader by АБВУ [Електронний ресурс] – Режим доступу: [https://ru.wikipedia.org/wiki/ABBY\\_FineReader](https://ru.wikipedia.org/wiki/ABBY_FineReader).
7. SimpleOCR system [Електронний ресурс] – Режим доступу: <http://freeanalogs.ru/SimpleOCR>.
8. Free-OCR system [Електронний ресурс] – Режим доступу: <http://www.free-ocr.com/>.
9. OCRFeeder system source code [Електронний ресурс] – Режим доступу: <https://github.com/GNOME/ocrfeeder>.
10. Tesseract-OCR source code [Електронний ресурс] – Режим доступу: <https://github.com/tesseract-ocr>.
11. Cloud Vision documentation [Електронний ресурс] – Режим доступу: <https://cloud.google.com/vision/docs>.

					ІАЛЦ.467100.003 ПЗ	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		57

12. Афонсенко А. В. Обзор методов распознавания структурированных символов [Электронный ресурс] / А. В. Афонсенко, А. И. Елизаров. – 2008. – Режим доступа до ресурсу: <http://old.tusur.ru/filearchive/reports-magazine/2008-2-1/83-88.pdf>.
13. Котович Н. В. Распознавание скелетных образов [Электронный ресурс] / Н. В. Котович, О. А. Славин. – 2003. – Режим доступа до ресурсу: <http://ocrai.narod.ru/skeletrecognize.html>.
14. How to use the tools provided to train Tesseract 4.00 [Электронный ресурс] – Режим доступа до ресурсу: <https://github.com/tesseract-ocr/tessdoc/blob/master/TrainingTesseract-4.00.md>.
15. Improving the quality of the output [Электронный ресурс] – Режим доступа до ресурсу: <https://tesseract-ocr.github.io/tessdoc/Technical-Documentation>.
16. Detect text in images [Электронный ресурс]. – Google, 2020. – Режим доступа: [https://cloud.google.com/vision/docs/ocr#vision\\_text\\_detection-python](https://cloud.google.com/vision/docs/ocr#vision_text_detection-python).
17. Document Text Tutorial [Электронный ресурс]. – Google, 2020. – Режим доступа: <https://cloud.google.com/vision/docs/fulltext-annotations>.
18. Data Usage FAQ [Электронный ресурс]. – Google, 2020. – Режим доступа: <https://cloud.google.com/vision/docs/data-usage>
19. Bots: An introduction for developers [Электронный ресурс]. – Telegram, 2020. – Режим доступа: <https://core.telegram.org/bots>.
20. Справочник по Bot API [Электронный ресурс]. – Telegram, 2020. – Режим доступа: <https://tlgrm.ru/docs/bots/api#authorizing-your-bot>
21. Python Telegram Bot's documentation [Электронный ресурс]. – Leandro Toledo, 2020. – Режим доступа: <https://python-telegram-bot.readthedocs.io/en/stable/>
22. Python Telegram Bot's Github [Электронный ресурс]. – Leandro Toledo, 2020. – Режим доступа: <https://github.com/python-telegram-bot/python-telegram-bot>

					ІАЛЦ.467100.003 ПЗ	Арк.
						58
Зм.	Арк.	№ докум.	Підпис	Дата		

23. Google. Setup the Vision API [Электронный ресурс] / Google – Режим доступа до ресурсу: <https://cloud.google.com/vision/docs/setup>.
24. Pycountry [Электронный ресурс]. – Christian Theune, 2019. – Режим доступа: <https://pypi.org/project/pycountry/>
25. The Heroku Platform [Электронный ресурс] – Режим доступа до ресурсу: <https://www.heroku.com/platform>.
26. Container Registry & Runtime (Docker Deploys) [Электронный ресурс] – Режим доступа до ресурсу: <https://devcenter.heroku.com/articles/container-registry-and-runtime>.

					<i>ІАЛЦ.467100.003 ПЗ</i>	Арк.
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>		59

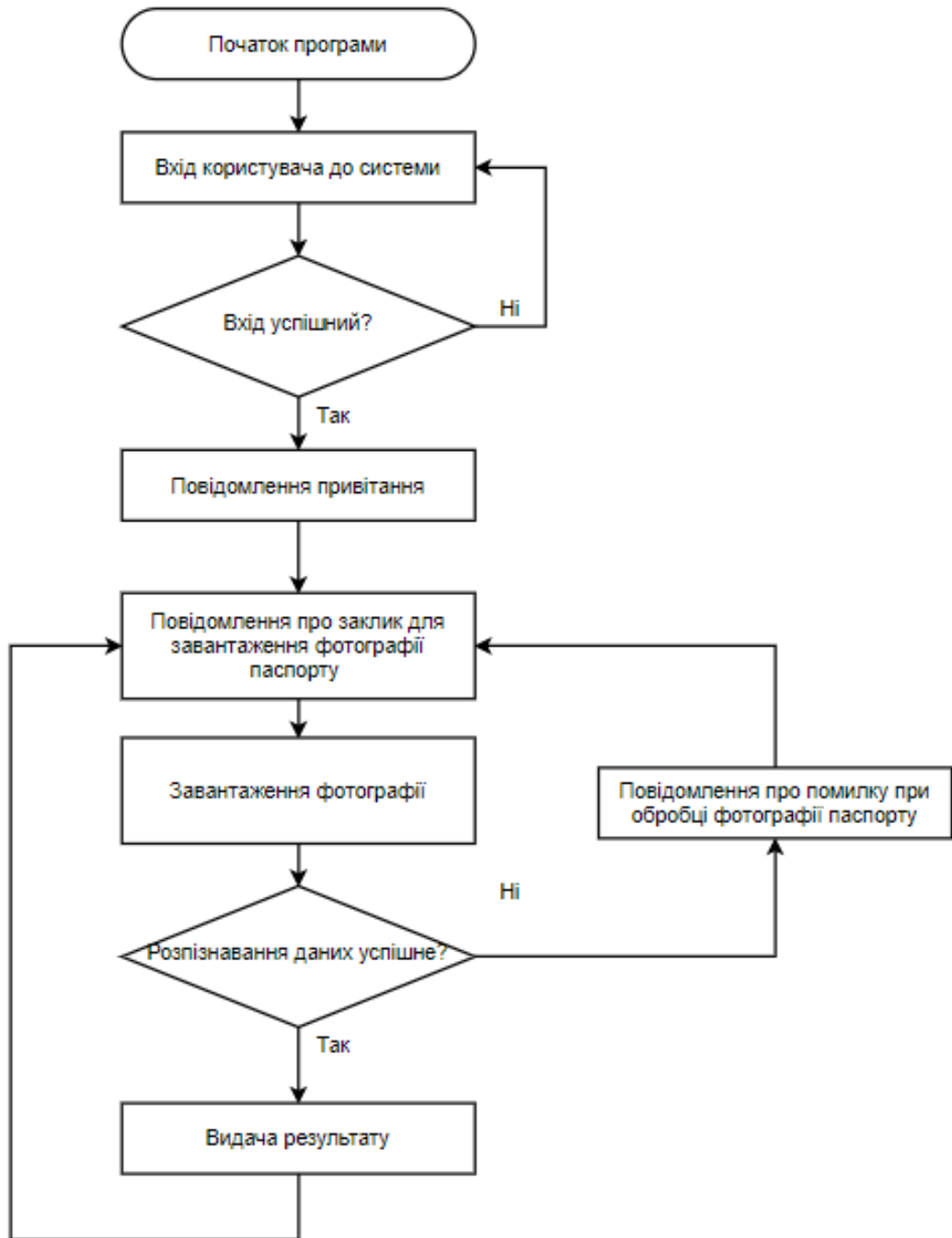
## **ДОДАТОК А**

«Система оптичного розпізнавання тексту на документах»

**Блок схема алгоритму**

Аркушів 1

Київ – 2020р.



					ІАЛЦ.467800.004 Д1		
Зм.	Арк.	№ докум.	Підпис	Дата			
Розробив		Бандурін В.Ю.			Лім.	Аркуш	Аркушів
Перевірив		Таран В.І.				1	1
Реценз.					НТУУ «КПІ», ФІОТ, ІО-64		
Н. Контр.		Сімоненко В.П.					
Затв.							
					Система розпізнавання тексту на документах Блок схема алгоритму		

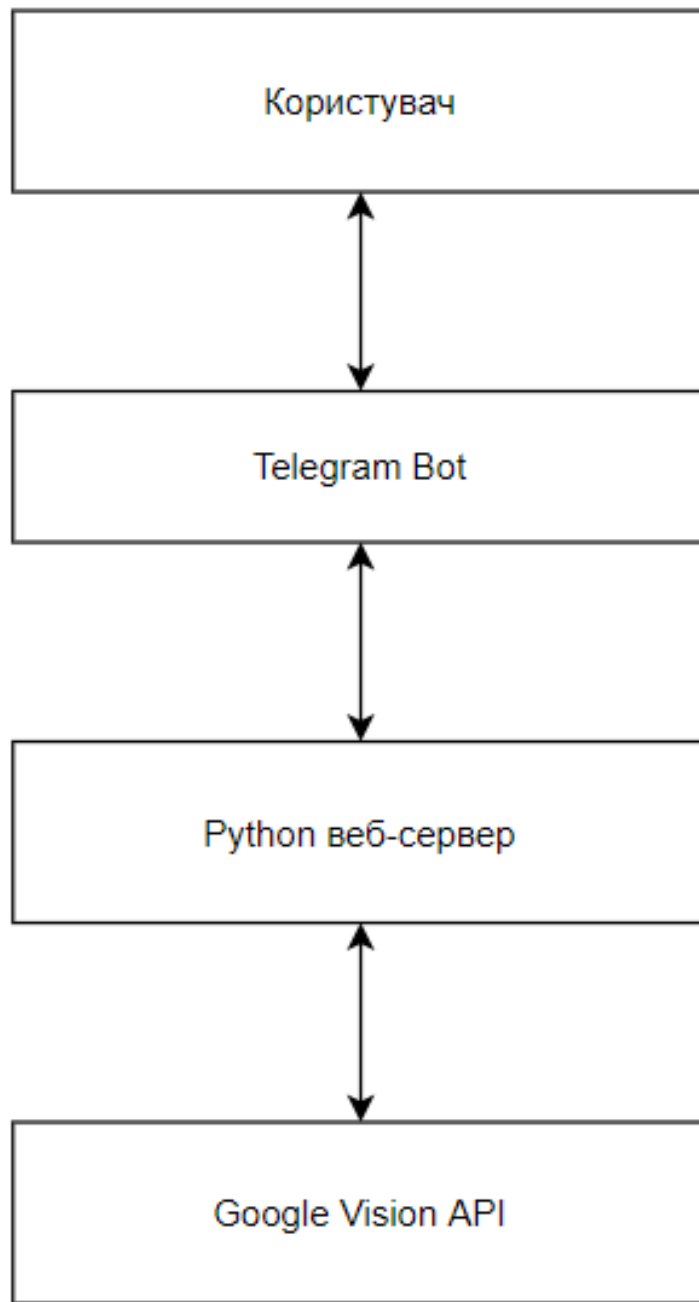
## **ДОДАТОК Б**

«Система оптичного розпізнавання тексту на документах»

**Схема роботи програми**

Аркушів 1

Київ – 2020р.



					ІАЛЦ.467800.004 Д2		
Зм.	Арк.	№ докум.	Підпис	Дата			
Розробив	Бандурін В.Ю.				Лім.	Аркуш	Аркушів
Перевірів	Таран В.І.					1	1
Реценз.					НТУУ «КПІ», ФІОТ, ІО-64		
Н. Контр.	Сімоненко В.П.						
Затв.							
					Система розпізнавання тексту на документах Схема роботи програми		

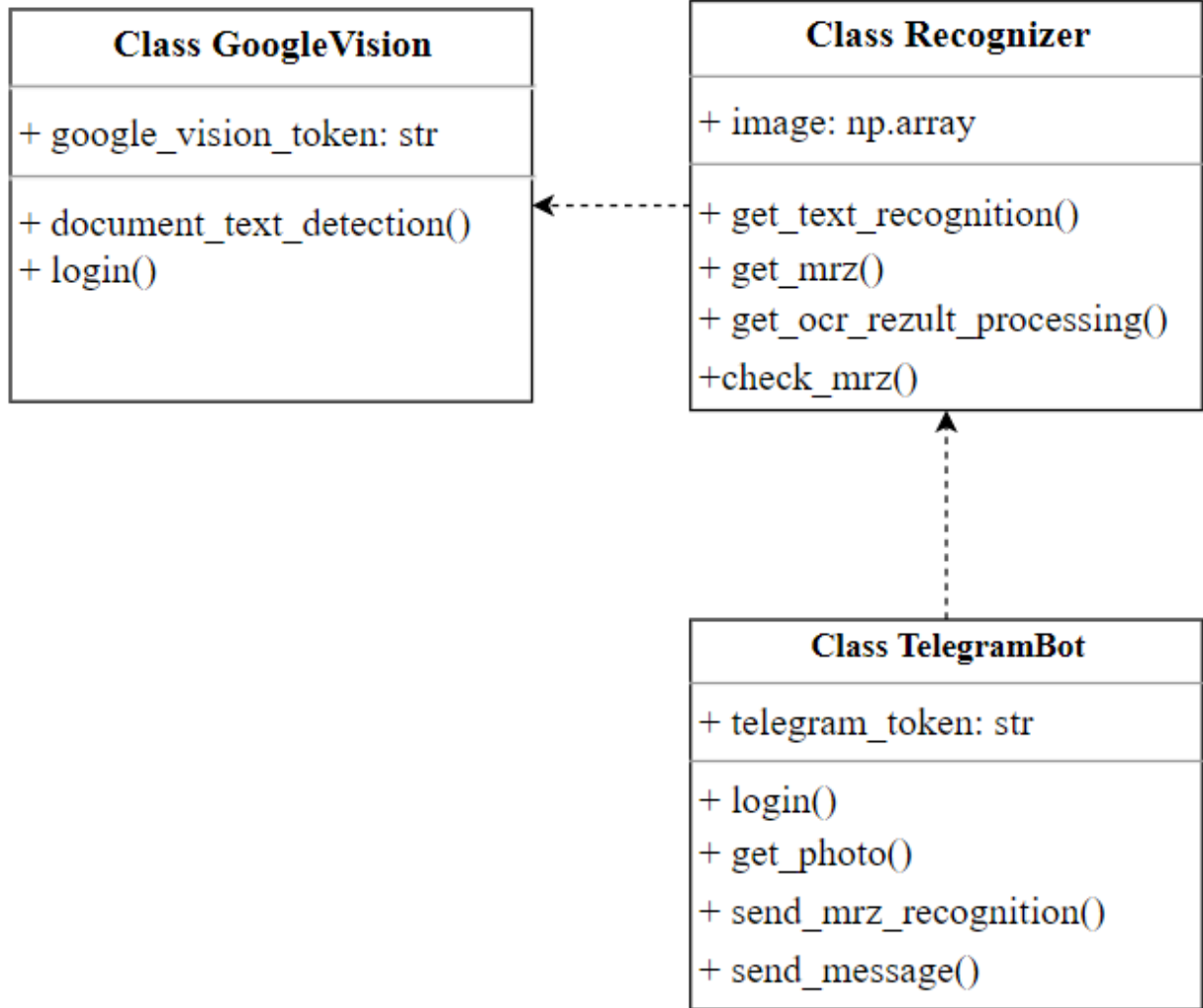
## **ДОДАТОК В**

«Система оптичного розпізнавання тексту на документах»

### **Діаграма класів**

Аркушів 1

Київ – 2020р.



					ІАЛЦ.467800.004 ДЗ			
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>	Система розпізнавання тексту на документах Діаграма класів	<i>Лім.</i>	<i>Аркуш</i>	<i>Аркушів</i>
Розробив	Бандурін В.Ю.						1	1
Перевірів	Таран В.І.							
Реценз.								
Н. Контр.	Сімоненко В.П.					НТУУ «КПІ», ФІОТ, ІО-64		
Затв.								

## ДОДАТОК Г

«Система оптичного розпізнавання тексту на документах»

Лістинг програми

Аркушів 32

Київ – 2020р.

```

import os
from google.cloud import vision
import settings
import re
os.environ["GOOGLE_APPLICATION_CREDENTIALS"] = settings.key

def recognize_document(img=None,
img_path=None):

    """
    Recognize document features in an
    image.
    :input img - numpy img or img_path
    - str path to img
    :return rez - json of recognized page
    """

    client =
vision.ImageAnnotatorClient()

    if img_path:
        with open(img_path, 'rb') as
image_file:
            content = image_file.read()
        else:
            content = img
        image =
vision.types.Image(content=content)

        response =
client.document_text_detection(image=
image)

        try:
            rez = response
            return rez

        except Exception as e:
            return {}

def get_country_code(text):
    country_code = re.findall('[0-9]+[A-
Z]{3}[0-9]+', text)
    if country_code:
        country_code = country_code[0]
        country_code = re.search('[A-

```

					ІАЛЦ.467800.004 Д4			
Зм.	Арк.	№ докум.	Підпис	Дата	Система розпізнавання тексту на документах Лістинг програми	Літ.	Аркуш	Аркушів
Розробив		Бандурін В.Ю.					1	32
Перевірив		Таран В.І.				НТУУ «КПІ», ФІОТ, ІО-64		
Реценз.								
Н. Контр.		Сімоненко В.П.						
Затв.								

```

Z]{3}', country_code).group(0)
else:
    country_code = None
return country_code
import common
import os
import json
import re
import pickle
import numpy as np
from enum import Enum
import time

import pycountry

class FeatureType(Enum):
    PAGE = 1
    BLOCK = 2
    PARA = 3
    WORD = 4
    SYMBOL = 5

class Lines():
    def __init__(self, mrz_block):
        bounding_boxes = []
        centres = []
        # print(mrz_block.bounding_box)

#
print(mrz_block.bounding_box.max(axis=0))
block_box =
mrz_block.bounding_box

max_x, max_y =
mrz_block.bounding_box.max(axis=0)
min_x, min_y =
mrz_block.bounding_box.min(axis=0)

center_block =
block_box.mean(axis=0)

for d in block_box:
    d = np.array(d)
    if d[0] == min_x:
        self.left_dot = d
    elif d[0] == max_x:
        self.right_dot = d
    elif d[1] == min_y:
        self.down_dot = d
    elif d[1] == max_y:
        self.top_dot = d

# print('dist',
self.get_dist(self.top_dot, self.left_dot))

for p in mrz_block.block:
    for w in p.paragraph:
        for s in w.word:
            box = s.bounding_box

```

```

        center = box.mean(axis=0)
        centres.append(center)
        #
    bounding_boxes.append(box)

```

```

def get_dist(self, d1, d2):
    return np.sqrt(((d1 - d2) ^
2).sum())

```

```

def get_rectangles(self):

```

```

    rectangle_1 = []
    rectangle_2 = []

```

```

class Symbol():

```

```

    def __init__(self, text, confidence,
bounding_box):
        self.text = text.upper()
        self.text =
common.translit(self.text)
        self.confidence = confidence
        self.bounding_box =
np.array([[bounding_box.vertices[0].x,
bounding_box.vertices[0].y],
[bounding_box.vertices[1].x,
bounding_box.vertices[1].y],

```

```

[bounding_box.vertices[2].x,
bounding_box.vertices[2].y],

```

```

[bounding_box.vertices[3].x,
bounding_box.vertices[3].y]])

```

```

class Word():

```

```

    def __init__(self, symbols_list,
bounding_box):

```

```

        self.word = symbols_list
        self.text = ".join([i.text for i in
self.word])

```

```

        self.confidence =
np.array([i.confidence for i in
self.word]).sum() / len(self.word)

```

```

        self.bounding_box =
np.array([[bounding_box.vertices[0].x,
bounding_box.vertices[0].y],

```

```

[bounding_box.vertices[1].x,
bounding_box.vertices[1].y],

```

```

[bounding_box.vertices[2].x,
bounding_box.vertices[2].y],

```

```

[bounding_box.vertices[3].x,
bounding_box.vertices[3].y]])

```

```

class Paragraph():

```

					IAJЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		

```

def __init__(self, word_list,
bounding_box):
    self.paragraph = word_list
    self.text = ".join([i.text for i in
self.paragraph])
    self.confidence =
np.array([i.confidence for i in
self.paragraph]).sum() /
len(self.paragraph)
    self.bounding_box =
np.array([[bounding_box.vertices[0].x,
bounding_box.vertices[0].y],
[bounding_box.vertices[1].x,
bounding_box.vertices[1].y],
[bounding_box.vertices[2].x,
bounding_box.vertices[2].y],
[bounding_box.vertices[3].x,
bounding_box.vertices[3].y]])

```

```

class Block():
    def __init__(self, paragraph_list,
bounding_box):
        self.block = paragraph_list
        self.text = ".join([i.text for i in
self.block])

```

```

self.confidence =
np.array([i.confidence for i in
self.block]).sum() / len(self.block)
    try:
        self.bounding_box =
np.array([[bounding_box.vertices[0].x,
bounding_box.vertices[0].y],
[bounding_box.vertices[1].x,
bounding_box.vertices[1].y],
[bounding_box.vertices[2].x,
bounding_box.vertices[2].y],
[bounding_box.vertices[3].x,
bounding_box.vertices[3].y]])
    except:
        self.bounding_box =
bounding_box

```

```

def get_words(self):
    words = []
    for p in self.block:
        for w in p.paragraph:
            words.append(w)
    return words

def get_symbols(self):
    symbols = []
    for p in self.block:
        for w in p.paragraph:

```

										IAЛЦ.467100.004 Д4
Зм.	Арк.	№ докум.	Підпис	Дата						Арк.

```

        for s in w.word:
            symbols.append(s)
    return symbols

def get_word_list(block_list,
mrz_blocks):
    word_list = []
    for b in block_list:
        ## print(len(b.get_words()),
len(mrz_block.get_words()))
        if b not in mrz_blocks:
            words = [w.text for w in
b.get_words()]
            word_list += words
        # else:
            #
print('=====
=====')
    return word_list

def save_json(j, file_name):
    dir_name = './recognized_json/'
    if file_name not in
os.listdir(dir_name):
        file_name = dir_name + file_name
        pickle.dump(j, file_name)
        # with open(file_name, 'w') as f:
        #     y = json.dumps(str(j))
        #     f.write(y)
    return True

def get_block_list(document):
    bounds = []

    block_list = []
    for page in
document.full_text_annotation.pages:
        for block in page.blocks:
            paragraph_list = []

            for paragraph in
block.paragraphs:
                word_list = []

                for word in paragraph.words:
                    symbols_list = []

                    for symbol in
word.symbols:
                        symbols_list.append(Symbol(symbol.te
xt, symbol.confidence,
symbol.bounding_box))

                bounds.append(symbol.bounding_box)

                bounds.append(word.bounding_box)
                w = Word(symbols_list,
word.bounding_box)
                word_list.append(w)
                p = Paragraph(word_list,
paragraph.bounding_box)

```

```

    paragraph_list.append(p)
    b = Block(paragraph_list,
block.bounding_box)

bounds.append(block.bounding_box)
    block_list.append(b)

return block_list

def get_mrz(block_list):
    blocks = []
    paragraph_list = []
    for b in block_list:
        text = b.text
        if '<' in text:
            blocks.append(b)
    if blocks:
        paragraph_list = []
        for b in blocks:
            paragraph_list += b.block

        b = Block(paragraph_list,
b.bounding_box)
        return b, blocks
    else:
        return []

```

```

class MRZ_block():

    def __init__(self, block_list,
word_list=None,
img='./images_24_06/2019-6-24
08_01_02.jpg'):

        mrz_block, blocks =
get_mrz(block_list)
        self.word_list =
get_word_list(block_list, blocks)
        self.mrz_block = mrz_block
        self.text = mrz_block.text

    for b in block_list:
        if b != mrz_block:
            self.not_mrz = b
            break

self.get_lines(img=img)
self.change_symbols()
self.check_country_numb()
self.check_passport_number()
self.check_birthday()
self.check_sex()
self.check_end_date()

def _find_similar(self, s1, s2):
    n = 0

```

					<i>IAJЦ.467100.004 Д4</i>	Арк.
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>		

```

if len(s1) <= len(s2):
    for i in range(len(s1)):
        if s1[i] == s2[i]:
            n += 1
    return n

def get_lines(self, plot=False,
img='../images_24_06/2019-6-24
08_01_02.jpg'):
    """
    Находит центры каждого
    символа. Записывает координаты
    каждого центра в centres.
    Словарь symbols_dict: key -
    center, value - Symbol
    Обучает линейный бинарный
    классификатор для
    классифицирования точек для
    каждой линии.
    (Линейная регрессия для всех
    центров, выше линии - один класс,
    ниже - другой)

    :return:
    np.array([first_line, second_line])
    """

text = self.mrz_block.text
text_cleaned = text.lstrip('<')

# print('text', text)

```

```

try:
    sep = re.findall('<{4,}',
text_cleaned)[0]
except:
    sep = '<<<<'

split_text = text.split(sep)
text_line_1, text_line_2 =
split_text[0], sep.join(split_text[1:])

symbols_list = []

for p in self.mrz_block.block:
    for w in p.paragraph:
        for s in w.word:
            symbols_list.append(s)

# center =
s.bounding_box.mean(axis=0)
# # print(text)
ind_1 = text.find(text_line_1),
text.find(text_line_1) + len(sep) +
len(text_line_1)
ind_2 = text.find(text_line_2),
text.find(text_line_2) +
len(text_line_2)
self.first_line =
symbols_list[ind_1[0]: ind_1[1]]
self.second_line =
symbols_list[ind_2[0]: ind_2[1]]

```

```

    return np.array([self.first_line,
self.second_line])

def change_symbols(self,
conf_1=0.3, conf_2=0.3):

    for s in self.first_line:
        if s.confidence <= 1:
            # print(s.text, s.confidence)
            if s.text in
common.always_change_1.keys():
                s.text =
common.always_change_1[s.text]

        for s in self.first_line:
            if s.confidence < conf_1:
                if s.text in
common.similar_letters_1.keys():
                    s.text =
common.similar_letters_1[s.text]

                n = self.first_line.index(s)
                if n != len(self.first_line) - 1:
                    if s.text == 'P' and
self.first_line[n+1].text == 'C' and
self.first_line[n-1].text != 'P' and n < 4:
                        #
print('=====C=====')
    )

```

```

        self.first_line[n+1].text =
'<'

    for s in self.second_line:
        if s.confidence < conf_2:
            # print(s.text, s.confidence)
            if s.text in
common.similar_letters_2.keys():
                s.text =
common.similar_letters_2[s.text]
                n = self.second_line.index(s)
                ls = '0'
                rs = '0'
                kr = "
ge = "
ur = "
                if n != len(self.second_line)-1
and n != 0:
                    ls = self.second_line[n-
1].text
                    rs =
self.second_line[n+1].text
                    ur = self.second_line[n-
1].text + self.second_line[n + 1].text
                    if n < len(self.second_line) - 3:
                        kr =
self.second_line[n+1].text +
self.second_line[n+2].text

                            if n > 2:

```

									IAJLЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата						

```

        ge = self.second_line[n-
2].text + self.second_line[n-1].text
        if s.text == 'S' and
self.second_line[n - 1].text == '<':
            s.text = '5'

        if s.text == 'O':
            s.text = '0'

        if s.text == '0' and (not
ls.isdigit() and not rs.isdigit()) and ls
not in '<O' and rs not in '<O':
            s.text = '0'

        if kr == 'KR' and (s.text == '0'
or s.text == 'O'):
            s.text = 'U'
        elif ge == 'GE' and (s.text ==
'0'):
            s.text = 'O'
        elif ur == 'UR' and (s.text ==
'<'):
            s.text = 'K'

        # print(".join([s.text for s in
self.first_line]))
        # print(".join([s.text for s in
self.second_line]))

def check_sex(self):
    """

```

Проверяет наличие М/Ф  
 Если нет, то ставит М/Ф в  
 нужную позицию, проверяя перед  
 ЭТИМ на проверке `check_end_date`,  
 куда входит М/Ф  
 :return:  
 {'is\_true': is\_true,  
 'sex': sex\_section}  
 """  
 self.ch\_sex = False  
 self.sex = None  
 second\_line = ".join([s.text for s in  
self.second\_line])

```

def is_true_sex(second_line):
    sex_section = re.findall('[0-
9]{6,}[A-Z][0-9]{3,}', second_line)
    ## print('sex_section',
sex_section)
    if len(sex_section) > 0:
        sex_section = sex_section[0]
    else:
        return {'is_true': False,
                'sex': ""}
    sex_section = [i for i in
sex_section if not i.isdecimal()][0]
    is_true = sex_section in
'MMMF'
    return {'is_true': is_true,
            'sex': sex_section}

```

```

d = is_true_sex(second_line)                sex_section = d['sex']
is_true = d['is_true']                       self.second_line[n].text =
sex_section = d['sex']                       'F'

if not is_true:                               self.ch_sex = is_true
    n =                                       self.sex = sex_section
second_line.find(self.country_code) +       return {'is_true': is_true,
10                                          'sex': sex_section}

    second_line = second_line[:n] +
'M' + second_line[n+1:]

    ch =
self.check_end_date(second_line=second_line)

if ch['is_true']:
    d = is_true_sex(second_line)
    is_true = d['is_true']
    sex_section = d['sex']
    self.second_line[n].text = 'M'

else:
    second_line =
second_line[:n] + 'F' + second_line[n +
1:]

    ch =
self.check_end_date(second_line=second_line)

if ch['is_true']:
    d =
is_true_sex(second_line)

    is_true = d['is_true']

def get_check_number(self, text):
    sum = 0
    for i in range(len(text)):
        x_dict = {0: 7,
                  1: 3,
                  2: 1}
        ind = i % 3
        x = x_dict[ind]
        n = text[i]
        if n.isdigit():
            sum += int(n) * x
        else:
            sum +=
common.later_dict[n] * x
    return sum

def check_passport_number(self):
    self.ch_number = False
    self.pass_number = None

    second_line = ".join([s.text for s in
self.second_line])

```

						<i>ІАЛЦ.467100.004 Д4</i>	<i>Арк.</i>
<i>Зм.</i>	<i>Арк.</i>	<i>№ докум.</i>	<i>Підпис</i>	<i>Дата</i>			

```
country_code = self.country_code
if not country_code:
    country_code = '[A-Z]{3}'

    reg_ex = '^\\w+<{0,}[0-9]{0,1}' +
country_code + '[0-9]{7}'
    number_section =
re.findall(reg_ex, second_line)

    if len(number_section) > 0:
        number_section[0] =
''.join(number_section[0].split('<'))
        number_section =
number_section[0][:-10]

        if None:
            ind_check_number =
number_section.find(country_code) - 1
            check_number =
number_section[ind_check_number]
            if check_number.isdecimal():
                check_number =
int(check_number)
            elif check_number in
common.leter2digit.keys():
                check_number =
int(common.leter2digit[check_number]
)

        else:
            if number_section[-1].isdecimal():
                check_number =
int(number_section[-1])
            elif number_section[-1] in
common.leter2digit.keys():
                check_number =
int(common.leter2digit[number_section[-1]])
            else:
                return {'is_true': False,
                        'text': None,
                        'check_number':
None}
            else:
                return {'is_true': False,
                        'text': None,
                        'check_number': None}
            self.pass_number =
number_section[:-1]
            if self.country_code == 'UKR':
                if len(self.pass_number) == 9:
                    self.pass_number =
self.pass_number[:-1]

else:
```

										IAJLЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата							

```

n =
second_line.find(self.pass_number) +
len(self.pass_number)
    self.second_line[n].text = '<'

sum =
self.get_check_number(self.pass_number)

is_true = sum % 10 ==
check_number

if not is_true:
    for i in
range(len(self.word_list)):
        w = self.word_list[i]
        w_2 = ""
        if i != len(self.word_list) - 1:
            w_2 = w + self.word_list[i
+ 1]
            n =
self._find_similar(self.pass_number, w)
            n_2 =
self._find_similar(self.pass_number,
w_2)

            if len(self.pass_number) ==
n:
                new_number =
str(self.pass_number) +
str(check_number)
                m =
self._find_similar(new_number, w)
                if m == len(new_number):
                    # print('equal', w)
                    self.pass_number =
new_number
                    sum =
self.get_check_number(self.pass_number)
                    check_number = sum %
10
                    is_true = True
                    # print(sum)
                    elif len(self.pass_number) *
0.6 <= n:
                        sum =
self.get_check_number(w)
                        if sum % 10 ==
check_number:
                            self.pass_number = w
                            is_true = True
                            # print('similar', w)
                            elif len(self.pass_number) *
0.6 <= n_2:
                                sum =
self.get_check_number(w_2)
                                if sum % 10 ==
check_number:
                                    self.pass_number = w_2
                                    is_true = True

```

```

else:
    w_2 = w_2.replace('O',
# country_code = re.findall('[0-
'0')          9]+[A-Z]{3}[0-9]+', text)
            text = ".join([s.text for s in
                self.second_line])
            sum =
            if sum % 10 ==
check_number:
            country_code = re.findall('[0-
                self.pass_number =
                9<]+[A-Z]{3}[0-9]{7}', text)
                # # print(country_code)
                if country_code:
                    w_2
                    is_true = True
                    sum =
                    country_code =
                    self.get_check_number(self.pass_number
                    country_code[0]
                    er)
                    country_code = re.search('[A-
                    Z]{3}', country_code).group(0)
                    country =
                    if check_number or
                    pycountry.countries.get(alpha_3=country
                    str(check_number) == '0':
                    code)
                    self.country_code =
                    sn = self.pass_number +
                    country_code
                    str(check_number)
                    if country:
                    else:
                    sn = self.pass_number
                    self.ch_number = is_true
                    self.ch_country = True
                    return {'is_true': is_true,
                    return {'is_true': True,
                    'text': self.pass_number,
                    'country': country.name,
                    'check_number':
                    'text': text}
                    check_number}
                    else:
                    def check_country_numb(self):
                    return {'is_true': False,
                    self.ch_country = False
                    self.country_code = None

```

						IAJLЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата			

```

        'country': None,
        'text': text}

    else:

        return {'is_true': False,
                'country': None,
                'text': text}

def check_birthdate(self):
    self.ch_birthdate = False

    second_line = ".join([s.text for s in
self.second_line])

    # second_line = ".join(zones[3:])
    country_code = self.country_code
    if not country_code:
        country_code = '[A-Z]{3}'

    reg_ex = country_code + '[0-
9]{7}'
    birthday_section =
re.findall(reg_ex, second_line)
    if birthday_section:
        birthday_section =
birthday_section[0]
    else:
        return {'is_true': False,
                'text': None,
                'check_number': None}
    check_number =
int(birthday_section[-1])
        text = birthday_section[3:-1]
        sum =
self.get_check_number(text)
        is_true = sum % 10 ==
check_number

        self.ch_birthdate = is_true
        self.birthday = text
        return {'is_true': is_true,
                'text': text,
                'check_number':
check_number}

def check_end_date(self,
second_line=None):
    self.ch_end_date = False

    if not second_line:
        second_line = ".join([s.text for s
in self.second_line])

    end_date_section = re.findall('[0-
9]{7}[A-Z]{1}[0-9]{7}', second_line)
    if end_date_section:
        end_date_section =
end_date_section[0]
    else:
        return {'is_true': False,
                'text': None,
                'check_number': None}

```

					ІАЛЦ.467100.004 Д4		Арк.	
Зм.	Арк.	№ докум.	Підпис	Дата				

```

        check_number =                                given_names = names[n +
int(end_date_section[-1])                            1:]
        text = end_date_section[8:-1]                else:
        sum =                                          return
self.get_check_number(text)
        is_true = sum % 10 ==                        tr_words = []
check_number                                          for w in self.word_list:
        self.ch_end_date = is_true
        self.end_date = text                        tr_words.append(common.name_transli
        return {'is_true': is_true,                 t(w))
                'text': text,
                'check_number':                     for name in given_names:
check_number}                                       for tr_w in tr_words:
                                                m =
                                                self._find_similar(name, tr_w)
                                                j = tr_words.index(tr_w)
                                                try:
                                                if name ==
                                                self.word_list[j][: -1]:
                                                given_names[given_names.index(name
                                                )] = self.word_list[j]
                                                except:
                                                pass
                                                try:
                                                if len(name) - 2 < m:
                                                if tr_w in
                                                self.word_list and tr_w !=
                                                self.word_list[j]:

```

```

given_names[given_names.index(name
)] = tr_w
        # print(tr_w)
    except Exception as e:
        print('error', e)
    try:
        t = tr_w.split('II')[0] +
'II'
        tw =
self.word_list[j].split('II')[0] + 'II'

        if name == t[:-1]:
            # print(t)

given_names[given_names.index(name
)] = t
            elif name == tw[:-1]:

given_names[given_names.index(name
)] = tw
            elif name[:-1] ==
tw[:-1]:

given_names[given_names.index(name
)] = tw

        except Exception as e:
            print('error', e)
        # print(given_names)

        for tr_w in tr_words:
            m =
self._find_similar(surname, tr_w)
            j = tr_words.index(tr_w)
            try:
                if surname ==
self.word_list[j][:-1]:
                    surname =
self.word_list[j]
            except:
                pass
            try:
                t = tr_w.split('II')[0] +
'II'
                if surname == t[:-1]:
                    # print(t)
                    surname = t
            except:
                print('error')

            if len(surname) - 2 < m:
                j = tr_words.index(tr_w)

                if tr_w in self.word_list
and tr_w != self.word_list[j]:
                    surname = tr_w
                    # print(tr_w)

```

```

        surname = surname_1 +
surname
        if len(names[n-1]) < 2:
            p =
'{:<<2s}'.format(names[n-1])
        else:
            p = "
        first_line = p + surname +
'<<' + '<'.join(given_names)
        first_line =
'{:<<44s}'.format(first_line)

        for i in
range(len(self.first_line)):
            self.first_line[i].text =
first_line[i]
            # print(first_line)
        except Exception as e:
            pass
            # print('error', e)

    def get_check(self):
        check_1 = self.ch_number
        check_2 = self.ch_country
        check_3 = self.ch_end_date
        check_4 = self.ch_birthday
        check_5 = self.ch_end_date
        check_6 = self.ch_sex

        ch_list = [check_1, check_2,
check_3, check_4, check_5, check_6]
        is_right = all(ch_list)
        ## print(check_1, check_2,
check_3, check_4, check_5, check_6)
        pass_ch =
self.get_check_number(self.pass_number) % 10

        # clear_first_line
self.clear_first_line()

        # self.split_names()
self.change_name()

self.format_names()

        second_line = ".join([s.text for s in
self.second_line])
        second_line = self.pass_number +
str(pass_ch) + self.country_code +
".join(second_line.split(self.country_code)[1:])

        first_line = ".join([s.text for s in
self.first_line])
        mrz_text = first_line +
second_line
        country_code = self.country_code

        return {'is_right': is_right,
                'mrz_text': mrz_text,
                'country_code':
country_code,

```

```

        'pass_num':
self.pass_number,
        'sex': self.sex,
        'birthday': self.birthday,
        'end_date': self.end_date}

def get_mrz(self):
    s = ".join([s.text for s in
self.first_line])
    s += ".join([s.text for s in
self.second_line])
    return s

def split_names(self):
    try:
        first_line = ".join([s.text for s in
self.first_line])
        names = ".join([s.text for s in
self.first_line]).split(self.country_code)
        # print(self.word_list)
        # print(names)
        if len(names) > 1:
            names = names[1]
            names = names.split('<')
            names = [i for i in names if i
and i != '<']
            # print(names)

            if len(names) == 1:

                k_split = names[0].split('K')

        comb = []
        comb_2 = []
        for i in range(len(k_split)):
            c = 'K'.join(k_split[0:1+i])
            comb.append(c)
            c = 'K'.join(k_split[-1 - i:])
            comb_2.append(c)
        # print(comb)
        change_ind = -1
        for s in comb:
            if s in self.word_list:
                change_ind =
first_line.find(s) + len(s)

        for s in comb_2:
            if s in self.word_list:
                change_ind =
first_line.find(s) - 1

        if change_ind != -1:
            self.first_line[change_ind].text = '<'

        # else:
        #     k_split =
names[0].split('C')
        #     comb = []
        #     comb_2 = []
        #     for i in
range(len(k_split)):

```

					IAJLЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		

```

# c = 'C'.join(k_split[0:1
+ i])
# comb.append(c)
# c = 'C'.join(k_split[-1-
i:])
# comb_2.append(c)
#
# print(comb)
#
# for s in comb:
#     if s in self.word_list:
#         change_ind =
first_line.find(s) + len(s)
#     if change_ind != -1:
#
self.first_line[change_ind].text = '<'
# else:
#     k_split =
names[0].split('S')
#     comb = []
#     comb_2 = []
#     for i in
range(len(k_split)):
#         c =
'S'.join(k_split[0:1 + i])
#         comb.append(c)
#         c = 'S'.join(k_split[-
1 - i:])
#         comb_2.append(c)
#         print(comb)
#

# for s in comb:
#     if s in
self.word_list:
#         change_ind =
first_line.find(s) + len(s)
#     if change_ind != -1:
#
self.first_line[change_ind].text = '<'
except Exception as e:
    print('split_names')
    print(e)
def format_names(self):
    try:
        first_line = ".join([s.text for s in
self.first_line])
        names = ".join([s.text for s in
self.first_line]).split('<<')
        names = [i for i in names if i
and i != '<']
        if len(names) >= 2:
            print(names)
        else:
            names = ".join([s.text for s in
self.first_line]).split('<')
            names = [i for i in names if i
and len(i) > 2]

```

```

    print(names)
    if len(names) > 1 and
len(names) == 2:
        n =
first_line.find(names[1])
            self.first_line =
self.first_line[:n] + [self.first_line[n-1]]
+ self.first_line[n:]
                # first_line = [s.text for s
in self.first_line]
                    # print(first_line)
except Exception as e:
    print('format_names')
    print(e)

def clear_first_line(self):
    try:
        if (len(self.first_line) +
len(self.second_line)) > 40:
            first_line = self.first_line

            text = ".join([s.text for s in
first_line])

            country_code =
self.country_code

            s = "[PV]+[A-Z0-9<`]{1}" +
country_code

            text_split = re.split(s, text)
            if len(text_split) > 1:
                n = text.find(text_split[0])
+ len(text_split[0])

```

```

    if n > 0:
        self.first_line =
self.first_line[n:]

    except Exception as e:
        print('clear_first_line')

    print(e)

def __str__(self):
    return self.get_mrz()

import os
from google.cloud import vision
import settings
import re

os.environ["GOOGLE_APPLICATION_CREDENTIALS"] = settings.key

def recognize_document(img=None,
img_path=None):
    """
    Recognize document features in an
    image.

    :input img - numpy img or img_path
    - str path to img

    :return rez - json of recognized page
    """

```

```

        country_code = re.search('[A-
client =
vision.ImageAnnotatorClient()
        Z]{3}', country_code).group(0)
else:
        country_code = None
return country_code

if img_path:
        with open(img_path, 'rb') as
image_file:
        content = image_file.read()
import json
else:
        content = img
later_dict = { '<': 0,
image =
        'A': 10,
vision.types.Image(content=content)
        'B':11,
        'C':12,
        'D':13,
        'E':14,
        'F':15,
        'G':16,
        'H':17,
        'I':18,
        'J':19,
        'K':20,
        'L':21,
        'M':22,
        'N':23,
        'O':24,
        'P':25,
        'Q':26,
        'R':27,
        'S':28,
        'T':29,
        response =
client.document_text_detection(image=
image)
try:
        rez = response
return rez
except Exception as e:
        return {}

def get_country_code(text):
        country_code = re.findall('[0-9]+[A-
Z]{3}[0-9]+', text)
if country_code:
        country_code = country_code[0]

```

						ІАЛЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата			

```
'U':30,                u'K': u'<',
'V':31,                u'4': u'<',
'W':32,                u'B': u'8',
'X':33,                u'E': u'F',
'Y':34,                u'S': u'5',
'Z':35}                u'5': u'S',
                        u'I': u'1',
```

```
similar_letters_1 = {
```

```
    u'K': u'<',
    u'4': u'<',
    u'C': u'<',
    u'E': u'F',
    u'O': u'D',
```

```
}
```

```
always_change_1 = {
```

```
    u'5': u'S',
    u'1': u'I',
    u'6': u'G',
    u'/': u'T',
    u'7': u'T',
    u'2': u'Z',
    u'0': u'O',
    u'4': u'<',
    u'$': u'S',
    u'-' : u'<',
    u'Ý': u'Y'
```

```
}
```

```
similar_letters_2 = {
```

```
    u'1': u'I',
    u'/' : u'7',
    u'T': u'7',
    u'7': u'T',
    u'O': u'0',
    u'6': u'G',
    u'-' : u'<',
    u'%': u'4',
    u'Ý': u'Y'
```

```
leter2digit = {u'O': u'0', u'T': u'7', u'I':
u'1', u'B': u'8', u'S': u'5'}
```

```
def name_translit(text):
```

```
    text = text.upper()
    capital_letters = {
        u'A': u'A',
        u'A': u'A',
        u'Б': u'B',
        u'В': u'V',
        u'В': u'V',
```

							ІАЛЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата				

u'Г': u'H',  
u'Г': u'G',  
u'Д': u'D',  
u'Е': u'E',  
u'Е': u'E',  
u'Ё': u'E',  
u'^Є': u'YE',  
u'Є': u'IE',  
u'Ж': u'ZH',  
u'З': u'Z',  
u'З': u'Z',  
u'И': u'Y',  
u'Т': u'Г',  
u'Т': u'Г',  
u'^І': u'YГ',  
u'І': u'Г',  
u'^Й': u'Y',  
u'Й': u'Г',  
u'Л': u'L',  
u'Н': u'N',  
u'Н': u'N',  
u'К': u'K',  
u'К': u'K',  
u'М': u'M',  
u'М': u'M',  
u'О': u'O',  
u'О': u'O',  
u'П': u'P',  
u'P': u'R',  
u'P': u'R',  
u'R': u'R',

u'C': u'S',  
u'C': u'S',  
u'T': u'T',  
u'T': u'T',  
u'У': u'U',  
u'У': u'U',  
u'Ф': u'F',  
u'Х': u'KH',  
u'Х': u'KH',  
u'Ц': u'TS',  
u'Ч': u'CH',  
u'Ш': u'SH',  
u'Щ': u'SHCH',  
u'^Ю': u'YU',  
u'Ю': u'IU',  
u'^Я': u'YA',  
u'Я': u'IA',  
# u'>': u'<',  
# u'«': u'<',  
# u'=: u'F',  
  
}  
translit\_string = ""  
  
s = '^' + text[0]  
if s in capital\_letters.keys():  
 translit\_string += capital\_letters[s]  
elif text[0] in capital\_letters.keys():  
 translit\_string +=  
capital\_letters[text[0]]

```

for i in text[1:]:
    u'Y': u'Y',
    u'Г': u'F',
    if i in capital_letters.keys():
        u'III': u'M',
        translit_string +=
    u'>': u'<',
    u'«': u'<',
    capital_letters[i]
    # else:
    u'=: u'F',
    # translit_string += i
}
return translit_string
translit_string = ""

def translit(string):
    """ This function works just fine """
    capital_letters = {
    u'H': u'H',
    u'З': u'3',
    u'A': u'A',
    u'B': u'B',
    u'B': u'B',
    u'E': u'E',
    u'Ё': u'E',
    u'И': u'I',
    u'ї': u'I',
    u'I': u'I',
    u'K': u'K',
    u'M': u'M',
    u'O': u'O',
    u'P': u'P',
    u'T': u'T',
    u'X': u'X',
    u'C': u'C',
    u'Ч': u'4',

    for i in string:
        if i in capital_letters.keys():
            translit_string +=
        capital_letters[i]
        else:
            translit_string += i

    return translit_string

def update_rezult(file_name, rezult):
    with open('rezult.json', 'r') as f:
        j = f.read()
        print(j)
    update = json.loads(j)
    update[file_name] = rezult
    update = json.dumps(update)
    with open('rezult.json', 'w') as f:
        f.write(update)
    return True

```

## Docker file

```
FROM ubuntu:latest
RUN apt-get update -y
RUN apt-get install -y python3-pip
python3-dev build-essential
COPY ./app
WORKDIR /app
RUN pip3 install -r requirements.txt
ENTRYPOINT ["python3"]
CMD ["app.py"]
```

```
from flask import Flask, jsonify,
send_file
from flask import request

import base64
import json
import os
import shutil
import time
import recognizer
import result_processing
from settings import host, database,
user, password, TOKEN
```

```
app = Flask(__name__)
DEBUG = True
```

```
try:
    import db_connector
    connector =
db_connector.DB_connector(host,
database, user, password)
except Exception as e:
    print(e)

@app.route('/api', methods=['POST'])
def get():
    data = request.get_json()
    data = json.loads(data)
    authorization =
request.headers['Authorization']
    img =
base64.b64decode(data["img"])
    n = len(os.listdir('./images')) + 1
    file_name = './images/' + str(n) +
'.jpg'
    with open(file_name, 'wb') as f:
        f.write(img)

    # img = data['img']
    if authorization == TOKEN:

        try:
            ex = str(None)
            block_list =
result_processing.get_block_list(docu
ment=rez)
```

						IAJLЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата			

```

    mrz =
result_processing.MRZ_block(block_li
st, img=f)
    # mrz_text = mrz.get_mrz()

    mrz_info = mrz.get_check()
    is_right =
str(int(mrz_info['is_right']))
    pass_num =
str(mrz_info['pass_num'])
    country_code =
str(mrz_info['country_code'])
    birthday = mrz_info['birthday']
    end_date =
mrz_info['end_date']
    mrz_text = mrz_info['mrz_text']
    sex = mrz_info['sex']

except Exception as ex:
    print(ex)
    mrz_text = ''
    is_right = '0'
    # result_json = {'mrz': '',
'is_right': '0', 'exception': ex}
    country_code = 'None'
    pass_num = 'None'
    birthday = 'None'
    end_date = 'None'
    sex = 'None'
    if is_right == '1':
        error = 'false'

else:
    error = 'true'
    # print(is_right, mrz_text,
country_code, pass_num)
    t = str(time.time())

connector.add_passport('heroku_v3_tes
t', is_right, mrz_text, country_code,
file_name, pass_num, t)
    # result_check =
is_mistake.is_mistake(text)

try:
    result_json = {'mrz': mrz_text,
                    # 'is_right': is_right,
                    #
'country_code':country_code,
                    #
'pass_num':pass_num,
                    # 'sex':sex,
                    # 'birthday': birthday,
                    'end_date': end_date,
                    # 'exception': ex,
                    'error': error}

except:
    result_json = {'mrz': mrz_text,
                    # 'is_right': is_right,
                    # 'country_code':
country_code,

```

3м.	Арк.	№ докум.	Підпис	Дата	ІАЛЦ.467100.004 Д4	Арк.
-----	------	----------	--------	------	--------------------	------

```

        # 'pass_num':
pass_num,
        # 'sex': sex,
        # 'birthday': birthday,
'end_date': end_date,
        # 'exception': 'n',
        'error': error}

    # file_name = str(n) + '.jpg'
    # update =
common.update_rezult(file_name,
str(rezult_json))
    # print(update)
else:
    rezult_json = {'Autorization error':
True}

    return jsonify(rezult_json)

@app.route('/', methods=['POST'])
def test():
    data = request.get_json()
    data = json.loads(data)

    rezult_json = {'success': True}
    return jsonify(rezult_json)

from telegram.ext import Updater,
CommandHandler, MessageHandler,
Filters
import settings
import os

import requests
import base64
import json
import re

# from database import *
DEBUG = True
TEST = True

class Bot():

    def __init__(self):
        self.start_bot()

    def reply_to_start_command(self,
bot, update):
        user = update.message.from_user
        user_id = user.id

        update.message.reply_text(
            "Привіт! За допомогою мене
ти зможеш розпізнати MRZ код
закордонного паспорту 😊\n"
            "Все що треба - це відправити
фото закордонного паспорту",)

        start_photo_id =
'AgADAgADcKoxG2b_cEvmMbRDy

```

						ІАЛЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата			



```

else:
    text = 'Не вдалося правильно
ропізнати паспорт. \nСпробуйте
використати фото кращої якості або
під іншим кутом нахилу. \nРезультат
распізнавання:\n'
    text = text + data['mrz']
    try:

        sep = re.search('<{3,>', text)[0]
        first = text.split(sep)[0] + sep
        second = text.split(sep)[1]
        text = first + '\n' + second
    except:
        pass

bot.sendMessage(chat_id=update.mess
age.chat_id, text=text)

# m =
bot.sendMessage(chat_id=362212345,
text=text)

def start_bot(self):
    print('bot started')

    if DEBUG == False:
        PORT =
int(os.environ.get('PORT', '8443'))

```

```

PORT =
int(os.environ.get('PORT', '5000'))

updater =
Updater(settings.TELEGRAM_API_K
EY)

dp = updater.dispatcher

dp.add_handler(CommandHandler("sta
rt", self.reply_to_start_command))

dp.add_handler(MessageHandler(Filter
s.photo, self.photo_handler))

updater.start_webhook(listen="0.0.0.0"
,

                                port=PORT,

url_path=settings.TELEGRAM_API_
KEY)

updater.bot.set_webhook("https://passp
ort-bot.herokuapp.com/" +
settings.TELEGRAM_API_KEY)

#
updater.bot.set_webhook("https://test-
pass-bot.herokuapp.com/" +
settings.TELEGRAM_API_KEY)

```

```
    updater.idle()

else:

    updater =
Updater(settings.TELEGRAM_API_K
EY_test)

    dp = updater.dispatcher

dp.add_handler(CommandHandler("sta
rt", self.reply_to_start_command))

dp.add_handler(MessageHandler(Filter
s.photo, self.photo_handler))

#
dp.add_handler(CommandHandler("hel
p", self.get_advice))

    updater.start_polling()
    updater.idle()

if __name__ == "__main__":
    bot = Bot()
```

					ІАЛЦ.467100.004 Д4	Арк.
Зм.	Арк.	№ докум.	Підпис	Дата		