

**ХЕРСОНСЬКИЙ НАЦІОНАЛЬНИЙ ТЕХНІЧНИЙ УНІВЕРСИТЕТ**  
(повне найменування вищого навчального закладу)  
**ФАКУЛЬТЕТ ІНФОРМАЦІЙНИХ ТЕХНОЛОГІЙ ТА ДИЗАЙНУ**  
(повне найменування інституту, назва факультету (відділення))  
**КАФЕДРА ПРОГРАМНИХ ЗАСОБІВ І ТЕХНОЛОГІЙ**  
(повна назва кафедри (предметної, циклової комісії))

## **Пояснювальна записка**

до кваліфікаційної роботи

**магістра**  
(освітній рівень)

на тему: «Розробка програмних інструментів розпізнавання тексту  
по зображеннях документів»

Виконав: студент 2 курсу, групи 5зІР2  
спеціальності  
121 - «Інженерія програмного забезпечення»  
(шифр і назва спеціальності)

Шерстюк Володимир Григорович  
(прізвище та ініціали)

Керівник д.т.н., професор Жарікова М.В.  
(прізвище та ініціали)

Рецензент \_\_\_\_\_  
(прізвище та ініціали)

Херсон - 2020

## АНОТАЦІЯ

Загальний обсяг роботи: 85 аркушів основного тексту, 14 ілюстрації, 6 таблиць. При підготовці використовувалася література з 30 різних джерел.

Робота складається із вступу, трьох розділів та двох додатків.

Об'єкт дослідження – процес розпізнавання та вилучення специфічних даних із документа по зображенню.

Предмет дослідження – засоби автоматизації процесу розпізнавання та вилучення специфічних даних із документа.

Мета дослідження – розробка програмної системи розпізнавання та вилучення специфічних даних із зображень документів та тестування системи із метою аналізу її ефективності.

Для досягнення мети дослідження поставлено і вирішено такі завдання:

- аналіз інструментів розпізнавання тексту;
- дослідження особливостей та вимог системи розпізнавання реквізитів документу по зображенню;
- проектування та розробка системи розпізнавання реквізитів документу по зображенню.

Наукова новизна полягає у розробці новітніх ефективних алгоритмів вилучення реквізитів документа із його зображення.

Практична цінність полягає в тому, що отримана система розпізнавання реквізитів документів по зображенню на основі відкритої бібліотеки Nicosoft OCR може бути впроваджена у будь-яку корпоративну систему електронного документообігу на комерційній основі.

Ключові слова: РОЗПІЗНАВАННЯ ТЕКСТУ, КОМП'ЮТЕРНИЙ ЗІР, ОПТИЧНЕ РОЗПІЗНАВАННЯ СИМВОЛІВ, СИСТЕМА РОЗПІЗНАВАННЯ РЕКВІЗИТІВ, РОЗПІЗНАВАННЯ ЗОБРАЖЕННЯ.

## ЗМІСТ

ПЕРЕЛІК ТЕРМІНІВ ТА СКОРОЧЕНЬ .....	7
ВСТУП .....	9
<b>РОЗДІЛ 1. ОГЛЯД ІСНУЮЧИХ РІШЕНЬ</b> <b>Ошибка! Закладка не определена.</b>	
1.1. Загальні відомості про задачу розпізнавання реквізитів та технології комп'ютерного зору .....	<b>Ошибка! Закладка не определена.</b>
1.2. Значення OCR у бізнесі .....	<b>Ошибка! Закладка не определена.</b>
1.3. Огляд найпопулярніших реалізацій OCR.....	<b>Ошибка! Закладка не определена.</b>
1.3.1. ABBYY FineReader .....	<b>Ошибка! Закладка не определена.</b>
1.3.2. Tesseract.....	<b>Ошибка! Закладка не определена.</b>
1.3.3. Asprise.....	<b>Ошибка! Закладка не определена.</b>
1.3.4. Nicomsoft.....	<b>Ошибка! Закладка не определена.</b>
Висновок до розділу 1.....	29
<b>РОЗДІЛ 2. АНАЛІЗ ЗАДАЧІ ТА ОГЛЯД ТЕХНОЛОГІЙ</b> <b>Ошибка! Закладка не определена.</b>	
2.1. Опис предметної області .....	<b>Ошибка! Закладка не определена.</b>
2.1.1. Повний опис процесу.....	<b>Ошибка! Закладка не определена.</b>
2.1.2. Основні задачі та склад системи ..	<b>Ошибка! Закладка не определена.</b>
2.2. Вибір технологій та їх обґрунтування .....	<b>Ошибка! Закладка не определена.</b>
2.2.1. Вибір платформи та мови програмування.....	<b>Ошибка! Закладка не определена.</b>
2.2.2. Інтерфейси, функції та об'єкти бібліотеки Nicomsoft OCR ....	<b>Ошибка! Закладка не определена.</b>
2.2.3. Вибір бібліотеки для проектування інтерфейсу .....	40
2.2.3. Вибір СУБД .....	41
2.3. Повний цикл роботи програми .....	<b>Ошибка! Закладка не определена.</b>
Висновок до розділу 2.....	<b>Ошибка! Закладка не определена.</b>

РОЗДІЛ 3. ПРОЕКТУВАННЯ ТА РОЗРОБКА СИСТЕМИ .....	<b>Ошибка! Закладка не определена.</b>
3.1. Визначення вимог для проектування.....	<b>Ошибка! Закладка не определена.</b>
3.2. Опис функціоналу інтерфейсного та автоматичного режимів ..	<b>Ошибка! Закладка не определена.</b>
3.3. Прецеденти .....	49
3.3.1. Прецеденти високого рівня.....	48
3.3.2. Створення шаблону .....	50
3.3.3. Робота програми у автоматичному режимі.....	55
3.4. Проектування бази даних.....	<b>Ошибка! Закладка не определена.</b>
3.5. Основні рішення з реалізації системи та її компонентів .....	59
3.5.1. Архітектура системи.....	59
3.5.2. Реалізація блоку розпізнавання тексту .....	62
3.5.3. Реалізація блоку пошуку шаблонів .....	63
3.5.4. Реалізація блоку вилучення реквізитів .....	63
3.6. Тестування системи .....	<b>Ошибка! Закладка не определена.</b>
Висновки до розділу 3 .....	<b>Ошибка! Закладка не определена.</b>
ЗАГАЛЬНІ ВИСНОВКИ .....	68
СПИСОК ВИКОРИСТАНИХ ДЖЕРЕЛ.....	71
ДОДАТОК А.....	74
ДОДАТОК Б .....	85

## ВСТУП

Сьогодні будь-яка сучасна організація чи компанія вводить електронний документообіг всередині компанії та з іншими компаніями. Кожен документ має певні реквізити та ключові слова, за якими його можна швидко знайти у базі даних. В залежності від типу документу, реквізити можуть бути різні, та представлені у різних типах даних. В будь-якому випадку, сьогодні майже у всіх організаціях, введенням реквізитів та даних по документам займається співробітник. Цей процес автоматизований хіба що, для компаній, в яких використовується один тип документу. Але коли типів таких документів багато - автоматизувати цей процес стає важче.

Розвиток нейронних мереж та систем комп'ютерного зору проходить саме зараз, але вже досяг певних результатів. Багато систем вміють досить точно розпізнавати об'єкти, що зображені на фотографіях, в тому числі і текст в документах.

Мета цієї магістерської кваліфікаційної роботи – вивчити можливість використання бібліотек із розпізнавання тексту для автоматизації розпізнавання їх реквізитів та розробка такої системи.

Існує велика кількість OCR бібліотек, потрібно було проаналізувати їх та порівняти. Потрібно було також розробити певний підхід, за яким система би могла «навчатись» точно викорінювати потрібні дані із документа.

Після цього вже можна було приступати до автоматизації бізнес-процесу із виявлення реквізитів та збереження їх у базі даних. Поміж таблиці із отриманими реквізитами у базі даних також потрібно було створити схему для збереження даних, потрібних саме для точного викорінення реквізитів.

Не існує єдиної, повноцінної системи, яка здатна використовувати OCR бібліотеки для автоматизації бізнес-процесу розпізнавання та викорінення реквізитів документів. Така система дозволить скоротити час на розпізнавання даних документів та їх збереження у базі даних. А висока точність розпізнавання також дає можливість проводити автоматично аналіз

документів. У світі, де намагаються автоматизувати всі робочі процеси, така система точно знайшла б своє місце.

Відсутності таких готових систем є пояснення. Зазвичай, у кожній компанії своя політика документообігу, та свої формати документів. Поміж розпізнавання, треба зробити аналіз документу та викорінити потрібні дані. При великій кількості різних типів документів це зробити досить важко. Тому єдиним вирішенням такої проблеми є створення у системі модуля, який би дозволяв «навчати» систему точно розпізнавати реквізити в кожному типі документу.

Звичайно, якість отриманих даних, а отже й якість всієї системи в цілому наряду залежить від точності розпізнавання символів. Тому був проведений аналіз найпопулярніших бібліотек із розпізнавання тексту, та підібраний найкращий із варіантів.

Вирішення цих ключових проблем та якісна технічна реалізація й приведе до створення системи із розпізнавання реквізитів документів по зображенню.

Система розроблялася як окремий модуль, який можна буде із легкістю впровадити у будь-яку систему із організації роботи компаній.