# Decognize: Prescription Digitization Using Knowledge Graphs

Project Team

Muhammad Sharjeel Akhtar   20P-0101
Mahad Ashraf                          20P-0563

Session 2020-2024

Supervised by

## Mr. Muhammad Sohaib Khan

**Department of Computer Science**

**National University of Computer and Emerging Sciences
Peshawar, Pakistan**

**June, 2024**

# Student's Declaration

We declare that this project titled "*Decognize: Prescription Digitization Using Knowledge Graphs*", submitted as requirement for the award of degree of Bachelors in Computer Science, does not contain any material previously submitted for a degree in any university; and that to the best of our knowledge, it does not contain any materials previously published or written by another person except where due reference is made in the text.

We understand that the management of Department of Computer Science, National University of Computer and Emerging Sciences, has a zero tolerance policy towards plagiarism. Therefore, We, as authors of the above-mentioned thesis, solemnly declare that no portion of our thesis has been plagiarized and any material used in the thesis from other sources is properly referenced.

We further understand that if we are found guilty of any form of plagiarism in the thesis work even after graduation, the University reserves the right to revoke our BS degree.

Muhammad Sharjeel Akhtar                    Signature: _____

Mahad Ashraf                                Signature: _____

_____

Verified by Plagiarism Cell Officer
Dated:

# Certificate of Approval

The Department of Computer Science, National University of Computer and Emerging Sciences, accepts this thesis titled *Decognize: Prescription Digitization Using Knowledge Graphs*, submitted by Muhammad Sharjeel Akhtar (20P-0101), and Mahad Ashraf (20P-0563), in its current form, and it is satisfying the dissertation requirements for the award of Bachelors Degree in Computer Science.

**Supervisor**

Mr. Muhammad Sohaib Khan                    Signature: _____

_____

Sir.Haroon Zafar

FYP Coordinator
National University of Computer and Emerging Sciences, Peshawar

_____

Dr. Noman Azam

HoD of Department of Computer Science
National University of Computer and Emerging Sciences

# Acknowledgements

Your acknowledgments here

Muhammad Sharjeel Akhtar

Mahad Ashraf

# Abstract

Doctors often write in incomprehensible handwriting, posing challenges for both the general public and pharmacists in understanding prescribed medications. Given the demanding nature of their work, doctors may not have the luxury to meticulously write prescriptions for the multitude of patients they attend to daily, resulting in illegible handwriting.To address this issue, "Decognize: Prescription Digitization Using Knowledge Graphs" tackles this problem by integrating Optical Character Recognition (OCR), deep learning and knowledge graphs. The system operates as an autonomous application, processing prescription images through OCR and image pre-processing for multiple languages. For handling illegible handwritten texts, advanced deep learning techniques i.e Transformers are employed to detect characters and interpret words. The system utilizes Natural Language Processing (NLP) knowledge embedded in knowledge graphs, establishing meaningful relationships between medicines based on their intended purposes. The key results demonstrate a notable improvement in the accuracy and comprehension of medical prescriptions. The system successfully recognizes medication names, reducing the likelihood of misspellings caused by illegibility. "Decognize" represents a transformative solution that streamlines the prescription digitization process, providing a user-friendly application for healthcare professionals and patients.

# Contents

# List of Figures

# List of Tables

# Chapter 1

# Introduction

- Develop NLP-based system for accurate prescription transcription.

- Inefficient healthcare data management for prescriptions.

- Errors in healthcare due to traditional OCR systems.

- This approach combines OCR and NLP, with BERN adding specialized biomedical knowledge and context from biomedical knowledge graph

- Improve scanned image quality and transcribe medical handwriting accurately using advanced models.

## 1.1 Purpose

In response to the evolving landscape of healthcare and the growing demand for efficient, data-driven solutions, our project represents a significant step forward in addressing a wide array of critical needs. From enhanced medical record digitization and improved healthcare decision-making to creating new research opportunities and enabling significant time and efficiency savings, our innovative approach is poised to provide a multifaceted solution. Moreover, our project is not merely a technological endeavor; it is a commitment to enhancing accessibility, fostering competitive recognition, and making a tangible contribution to local healthcare. This introduction sets the stage for a project

that not only leverages cutting-edge technology but also places the wellbeing of patients and the efficiency of healthcare practitioners at the forefront, ultimately leading to an improved and more responsive healthcare ecosystem.

### 1.1.1 Motivation

The motivation behind this project is rooted in a profound recognition of the critical healthcare challenges that persist in our current healthcare system. By focusing on aspects such as improving data management, preventing medication errors, and ensuring swift access to precise patient information, our project aims to address fundamental issues that impact the quality of care. Furthermore, by facilitating informed decision-making and reducing the burdens of manual data entry, we strive to enhance the efficiency and effectiveness of healthcare delivery. The project also seeks to bridge the gap between healthcare and technology research, opening up new possibilities for innovation and advancement in the field. Handwritten prescription handling is a particularly challenging aspect of healthcare documentation, and our approach leverages Knowledge Graph technology to tackle this issue head-on. Ultimately, the overarching goal is to significantly enhance patient safety and healthcare quality, driving a positive transformation in the way healthcare data is managed, leading to safer, more efficient, and higher-quality healthcare outcomes.

# Chapter 2

# Review of Literature

## 2.1 OCR Enhancement Using Tesseract and OpenCV (2023)

This paper explores the enhancement of OCR using Tesseract and OpenCV with a focus on medical handwriting recognition. The goal is to improve data management in healthcare records. The methodology integrates image pre-processing and text detection. Results indicate satisfactory accuracy in well-preprocessed images but struggles with complex backgrounds and artifacts.

### 2.1.1 Conclusion:

While effective in recognizing characters, Tesseract's accuracy is hindered by poor image quality and the need for meticulous preprocessing, which may limit its applicability in complex cases.

## 2.2 OCR with TensorFlow and Custom ResNet Model (2021)

This research focuses on enhancing OCR using TensorFlow and a custom ResNet model, aiming to improve healthcare data management and patient safety. This research focuses

on enhancing OCR using TensorFlow and a custom ResNet model, aiming to improve healthcare data management and patient safety. The methodology enhances model robustness with data augmentation and custom architecture. The results demonstrate robustness and suitability for the specified task.

### 2.2.1 Conclusion:

In conclusion, our research successfully improved OCR using TensorFlow and a custom ResNet model. The implemented methodology, including data augmentation and a specialized architecture, demonstrated robustness and suitability for healthcare data management, thus contributing to enhanced patient safety.

## 2.3 Biomedical Knowledge Graph Construction with NLP (2021)

This paper introduces the construction of a biomedical knowledge graph with NLP, emphasizing text extraction from biomedical documents. The methodology involves OCR, BERN, and zero relation extraction. Successful establishment of a Neo4j knowledge graph is highlighted.

### 2.3.1 Conclusion:

While versatile in applications, limitations include potential inaccuracies in the zero-shot relation extractor, requiring expert validation.

## 2.4   Multilingual Handwritten Prescription Recognition (2022)

This research focuses on using CNNs, RNNs, and LSTMs to recognize and translate handwritten prescriptions in multiple languages. Machine learning techniques like CNNs and LSTMs are employed to achieve efficient recognition. Sensitivity to variations in handwriting styles is highlighted.

### 2.4.1   Conclusion:

While efficient, variations in handwriting styles and the quality and diversity of training data can impact performance.

## 2.5   Machine Learning Algorithms for Prescription Text (2022)

This study compares various machine learning algorithms for recognizing text on medical prescriptions. The methodology involves image scanning and CNN-based feature extraction, leading to successful recognition of medical prescriptions.

### 2.5.1   Conclusion:

While successful, the study acknowledges the need for further investigation into alternative machine learning algorithms for comprehensive comparison.

## 2.6   Refernces

- Pavithiran G, Sharan Padmanabhan, Nuvvuru Divya, Aswathy V, Irene Jerusha P, Chandar B. (2022).  "Doctor's Handwritten Prescription Recognition System In Multi-Language Using Deep Learning."

- Bratanic, T. (2021, October 25). Constructing a Biomedical Knowledge Graph with NLP. Towards Data Science. Retrieved from https://towardsdatascience.com/constructing-a-biomedical-knowledge-graph-with-nlp-9c07f13d275e

- Hassan, Esraa, Habiba Tarek, Mai Hazem, Shaza Bahnacy, Lobna Shaheen, and Walaa H.Elashmwai.(2021) "Medical prescription recognition using machine learning. In 2021 IEEE 11th Annual Computing and Communication Workshop and Conference (CCWC), pp. 0973-0979. IEEE

- Filip Zelic and Anuj Sable Areviewonon (2023) OCR with Tesseract OpenCV and Python.Nanonets.

| Sr. no | Year | Basic Idea | Methodologies | Results | Limitations |
|--------|------|------------|---------------|---------|-------------|
| [1] | 2023 | OCR with Open CV and tesseract | Implemented Tesseract OCR with Open CV in python. Focusing on image pre-processing for optimal results integrated text detection and recognition components | Achieved satisfactory OCR accuracy with well-preprocessed images. However, Tesseract struggled with complex backgrounds and artifacts, yielding suboptimal outputs. | Tesseracts accuracy is hindered by poor image quality. Requiring meticulous preprocessing . Challenges arise in handling artifacts handwriting and diverse languages . |
| [2] | 2021 | Optical Character Recognition Using TensorFlow | Implemented OCR with TensorFlow Enhanced model robustness with data augmentation technique. Implemented a custom ResNet architecture for OCR | These results showcase the effectiveness of the OCR model, particularly in accurately recognizing characters within the test set, demonstrating its robustness and suitability for the specified task. | Our Model can fail if the image is complex . E.g cursive writing images or images with Continous Characters Currently our model is trained only on digits and English language |
| [3] | 2021 | Construct a Bio Medical Knowledge Graph with NLP | Extracted text from biomedical document using OCR and applied BERN and utilized zero relation extractor. | Successfully established a Neo4j knowledge graph, showcasing versatility through demonstrated applications such as search engine, co-occurrence analysis and author expertise inspection. While emphasizing its utility for diverse biomedical machine learning applications. | Limitations include persistent NER challenges with BERN, potential inaccuracies in the zero shot relation extractor and the need for expert validation with external database enrichment reliant on data consistency |
| [4] | 2018 | Build a Handwritten Text Recognition System using TensorFlow | Implemented HTR using TensorFlow, with NN trained on IAM word-images, including CNN, RNN and CTC layers. Preprocessed data with resizing normalization and potential augmentation. Utilized RMSProp for training and explored enhancements like data augmentation, input size adjustments and decoding strategies. | Implemented successful HTR on IAM word -images, enabling flexible NN customization and identifying areas for accuracy improvements. | Limited Diversity due to reliance on IAM dataset Potential recognition errors especially for non-dictionary words CPU based training may be slower : GPU recommended |
| [5] | 2022 | Doctor Handwritten Prescription recognition system in multilanguage using deep learning | Implemented a system employing machine learning techniques such as CNNs,RNNs,LSTMs for recognizing and translating handwritten prescription notes in diverse language | Successful recognition and translation of handwritten prescriptions in various languages. Demonstrated the efficiency of CNNs, RNNs, and LSTMS in multilingual handwritten text processing. | Sensitivity to variations in handwriting styles. Reliance on quality and diversity of training data for optimal performance |
| [6] | 2022 | A Comparison of various Machine learning Algorithms for recognizing Text on Medical prescriptions | Proposed approach involves image scanning pre-processing and CNN-based feature extraction for recognizing handwritten medical prescriptions. Results are compared with drug name database using OCR for medicinal name identification | Successful implementation of CNN-based recognition for medical prescription. Need for further investigation into alternative machine learning algorithms for comprehensive comparision. | Limited Exploration of alternative machine learning algorithms Identification challenges with low accuracy medical names in OCR |
| [7] | 2020 | Online Cursive Handwritten Medical Words Recognition System | Implemented an online cursive handwritten medical word recognition system using a bidirectional LSTM network. Employed data augmentation techniques to enhance recognition efficiency. | Successful Utilization of bidirectional LSTM for cursive medical word recognition Recognition efficiency improvements achieved Through data augmentation | The system is restricted to providing output only for the trained data Inability to generate output for the new unseen data due to lack of adaptability |
| [8] | 2021 | Medical Prescription Recognition Using Machine Learning | Developed a Medical Prescription Recognition System employing image processing techniques and machine learning algorithms to identify handwritten medicine names from prescription note images. | Successful integration of image processing and machine learning for medical prescription recognition . Acknowledged limitations include reliance on small dataset and lower accuracy | Limited dataset usage in the system The system exhibits low accuracy levels |
| [9] | 2021 | Medical Prescription Identification Solution | Implemented a Medical Prescription Identification Solution employing a neural network for character recognition and knowledge-based maching for accurate results. | Utilized neural network approach and knowledge based matching for effective prescription identification | Restricted to reading only one line at a time |

Table 2.1: Literature Review
Column 1 defines the citation. Column 2 defines the year. Column 3 defines the basic idea. Column 4 defines the methodology applied. Column 5 defines the results obtained. Column 6 defines the limitations.

# Chapter 3

# Project Vision

## 3.1 Problem Statement

### 3.1.1 Problem

Inefficient healthcare data management for prescriptions.

### 3.1.2 Challenge

Illegible handwriting , medical jargon and Knowledge Graph

### 3.1.3 Consequences

Errors in healthcare due to traditional OCR systems.

### 3.1.4 Goal

Develop NLP-based system for accurate prescription transcription

## 3.2   Business Opportunity

### 3.2.1   Efficient Data Digitization:

Streamline the conversion of handwritten medical notes into digital records, saving time and reducing errors.

### 3.2.2   Enhanced Decision-Making:

Improve data accuracy for quicker, informed healthcare decisions, boosting patient care.

### 3.2.3   Research and Efficiency:

Enable advanced research and analytics with biomedical knowledge graphs, driving medical innovation.

### 3.2.4   Cost Savings:

Automate manual data entry, significantly cutting operational costs for healthcare institutions.

### 3.2.5   Competitive Edge:

Adopting this technology positions healthcare providers as industry leaders.

### 3.2.6   Scalable Deployment:

Tailor the solution for diverse healthcare institutions, from clinics to hospitals.

### 3.2.7  Global Reach:

Multilingual support ensures accessibility in international markets.

## 3.3   Objectives

- To reduce error percentage in prescriptions readability.

- To create an improved OCR system which could later on deployed on other real-life-domains as well.

- To allow user to save and access their prescription data conveniently.

## 3.4   Project Scope

- Global OCR in Healthcare: Booming market, 15.4

- Recent Projects: Automated doctor prescription by Nano Net Technologies Inc and Neurodata Group.

- OCR in Healthcare in Pakistan: Active research by Seerat Rani, Abd Ur Rehman, Beenish Yousaf, Hafiz Tayyab Rauf, Emad Abouel Nasr, and Seifedine Kadry.

- Summary: OCR enhancing healthcare in Pakistan through innovation and integration.

## 3.5   Constraints

### 3.5.1  Data Privacy and Compliance:

The project must adhere to stringent healthcare data privacy regulations, such as HIPAA and GDPR, which can limit data sharing and access.

### 3.5.2 Data Quality Variability:

Variability in the quality of handwritten medical notes and scanned images may impact the accuracy of OCR and NLP, especially in cases of poor handwriting or low-resolution scans.

### 3.5.3 Multilingual and Handwriting Variability:

Healthcare documents may be written in multiple languages and various handwriting styles, posing a challenge for accurate recognition and translation.

### 3.5.4 Hardware and Connectivity:

Accessibility to compatible hardware and a reliable internet connection may be a constraint, especially for smaller healthcare providers with limited resources.

### 3.5.5 Initial Investment:

Implementing and customizing the system may require a substantial initial financial investment, which can be a constraint for budget-constrained healthcare institutions.

### 3.5.6 User Training and Adoption:

Healthcare professionals and staff may require training to effectively use the technology, and resistance to adopting new systems can be a constraint.

### 3.5.7 Integration with Existing Systems:

Seamless integration with existing Electronic Health Record (EHR) systems, with their unique standards and formats, can pose integration challenges.

### 3.5.8   Maintenance and Updates:

Ongoing maintenance, updates, and IT support are crucial for system reliability, and the associated costs and resource requirements can be constraints.

### 3.5.9   Scalability:

Ensuring that the system can scale to accommodate growing data volumes and user loads without performance degradation can be a challenge.

### 3.5.10   Vendor Lock-In:

Dependence on a specific vendor for the technology may limit flexibility and pose long-term constraints.

### 3.5.11   Data Sovereignty:

Geopolitical factors and data sovereignty laws may restrict the storage location and cross-border transfer of healthcare data.

### 3.5.12   Ethical and Legal Considerations:

The application of advanced technologies in healthcare data management may raise ethical and legal concerns, potentially constraining certain aspects of the project's implementation.

### 3.5.13   Accessibility in Remote Areas:

Healthcare institutions in remote or underserved areas may face challenges in accessing the technology due to infrastructure limitations.

## 3.6 Stakeholders Description

### 3.6.1 Healthcare Providers:

Healthcare institutions, including hospitals, clinics, and private practices, are primary stakeholders. They are interested in efficient data management, enhanced patient care, and cost savings.

### 3.6.2 Patients:

Patients are indirect stakeholders as they benefit from improved data accuracy, which contributes to better healthcare decision-making, reduced medication errors, and enhanced quality of care.

### 3.6.3 Healthcare Professionals:

Physicians, nurses, and other healthcare staff are critical stakeholders. They use the system to access and manage patient data, impacting their daily workflows and decision-making.

### 3.6.4 Health IT Vendors:

Companies providing healthcare information technology solutions, EHR systems, and data management software are stakeholders, as the project may complement or compete with their offerings.

### 3.6.5 Pharmaceutical and Biotech Companies:

Stakeholders in this sector can benefit from the biomedical knowledge graph for drug research and development.

### 3.6.6   Telemedicine Platforms:

Telemedicine providers are stakeholders interested in enhancing data management capabilities and accessibility in remote healthcare services.

### 3.6.7   Regulatory Authorities:

Regulatory bodies like the FDA, HIPAA, and GDPR play a crucial role in shaping the project's compliance and data privacy aspects.

### 3.6.8   Data Security Experts:

Cybersecurity and data privacy experts are stakeholders in ensuring the security and privacy of patient data within the project.

### 3.6.9   Research Institutions:

Academic and research institutions may benefit from the project's data for medical research and analytics, making them indirect stakeholders.

### 3.6.10   Investors and Funders:

Individuals or organizations investing in the project or providing funding are stakeholders with an interest in the project's financial success.

### 3.6.11   Technology Providers:

Suppliers of hardware, software, and infrastructure required for the project are stakeholders.

### 3.6.12    Data Scientists and Analysts:

Professionals with expertise in data analysis and NLP are essential stakeholders in developing, maintaining, and improving the project's capabilities.

### 3.6.13    Ethical and Legal Advisors:

Experts in healthcare ethics and legal matters are stakeholders, ensuring that the project complies with ethical and legal standards.

### 3.6.14    Patients' Advocacy Groups:

Groups advocating for patient rights and privacy are indirect stakeholders, with an interest in how the project impacts patient data.

### 3.6.15    Consulting and Integration Firms:

Firms providing consulting and integration services for the project's implementation are stakeholders.

# Chapter 4

# Software Requirements Specifications

## 4.1 List of Features

- **Optical Character Recognition (OCR):**

    – Accurate recognition of handwritten medical notes.

    – Image preprocessing for optimal OCR results.

    – Multilingual support for diverse patient populations.

- **Natural Language Processing (NLP):**

    – Extraction and categorization of medical entities (e.g., drugs, diseases, procedures).

    – Text disambiguation for improved readability.

    – Multilingual NLP capabilities.

- **Biomedical Knowledge Graph:**

    – Construction of a structured graph using standardized vocabularies.

    – Linking medical entities to enrich data.

    – Facilitation of data retrieval for research and analytics.

- **Integration with Electronic Health Records (EHR):**

- – Seamless integration with existing EHR systems.

- – Compatibility with various EHR standards and formats.

- **Data Privacy and Security:**

  - – Compliance with healthcare data privacy regulations (e.g., HIPAA, GDPR).

  - – Encryption and secure storage of patient data.

  - – Access control and audit trails.

- **Scalability:**

  - – Ability to scale with growing data volumes and user loads.

  - – Support for healthcare institutions of all sizes, from clinics to large hospitals.

- **User Training and Support:**

  - – Training resources and documentation for healthcare professionals.

  - – Ongoing technical support and assistance.

- **Multilingual Support:**

  - – Recognition and translation of multiple languages.

  - – Accommodation of diverse patient populations.

- **Data Quality Improvement:**

  - – Image enhancement to improve OCR accuracy.

  - – Data validation and correction mechanisms.

- **Efficient Data Retrieval:**

  - – Fast and accurate data retrieval for healthcare professionals.

  - – Advanced search capabilities for research and decision-making.

- **Cost and Time Savings:**

  - – Automation of manual data entry processes.

  – Reduction of operational costs for healthcare institutions.

• **Comprehensive Reporting:**

  – Generation of detailed reports for data analysis.

  – Customizable reporting options.

• **Regulatory Compliance Tools:**

  – Tools for ensuring compliance with healthcare regulations.

  – Regular updates to maintain compliance with evolving laws.

• **Adaptive Learning:**

  – Machine learning algorithms that adapt and improve over time.

  – Continuous optimization for better recognition accuracy.

• **Customization and Integration:**

  – Tailoring the system to suit the unique needs of each healthcare institution.

  – APIs for seamless integration with other healthcare systems.

• **Cross-Platform Accessibility:**

  – Accessible via web browsers and mobile devices.

  – Cloud-based solutions for remote access.

• **Research and Analytics Tools:**

  – Access to a rich biomedical knowledge graph for research and innovation.

  – Analytics features for healthcare data-driven decision-making.

## 4.2 Functional Requirements

• The system shall accurately recognize handwritten medical notes using OCR.

• It shall perform image preprocessing to optimize OCR results.

- OCR shall support multiple languages for diverse patient populations.

- The system shall extract and categorize medical entities, including drugs, diseases, and procedures using NLP.

- It shall disambiguate text for improved readability.

- NLP capabilities shall support multiple languages.

- The system shall construct a structured biomedical knowledge graph using standardized vocabularies.

- It shall link medical entities to enrich data.

- The system shall facilitate data retrieval for research and analytics.

- It shall seamlessly integrate with existing Electronic Health Records (EHR) systems.

- The system shall be compatible with various EHR standards and formats.

- It shall comply with healthcare data privacy regulations, such as HIPAA and GDPR.

- The system shall encrypt and securely store patient data.

- Access control and audit trails shall be implemented.

- The system shall be scalable to accommodate growing data volumes and user loads.

- It shall support healthcare institutions of all sizes, from clinics to large hospitals.

- The system shall provide training resources and documentation for healthcare professionals.

- Ongoing technical support and assistance shall be available.

- The system shall recognize and translate multiple languages.

- It shall accommodate diverse patient populations.

- The system shall enhance images to improve OCR accuracy.

- It shall include data validation and correction mechanisms.

- The system shall provide fast and accurate data retrieval for healthcare professionals.

- It shall offer advanced search capabilities for research and decision-making.

- The system shall automate manual data entry processes.

- It shall reduce operational costs for healthcare institutions.

- The system shall generate detailed reports for data analysis.

- It shall offer customizable reporting options.

- The system shall provide tools for ensuring compliance with healthcare regulations.

- It shall include regular updates to maintain compliance with evolving laws.

- The system shall use machine learning algorithms that adapt and improve over time.

- Continuous optimization shall be implemented for better recognition accuracy.

- The system shall allow tailoring to suit the unique needs of each healthcare institution.

- It shall offer APIs for seamless integration with other healthcare systems.

- The system shall be accessible via web browsers and mobile devices.

- It shall offer cloud-based solutions for remote access.

- It shall grant access to a rich biomedical knowledge graph for research and innovation.

- It shall include analytics features for healthcare data-driven decision-making.

21

## 4.3   Quality Attributes

- **Accuracy:** The system must ensure high accuracy in recognizing handwritten medical notes and categorizing medical entities.

- **Reliability:** Healthcare professionals must rely on the system for accurate data recognition and retrieval.

- **Security:** Robust data security measures must be in place to protect patient data and comply with privacy regulations.

- **Usability:** Healthcare professionals should find the system user-friendly and easy to navigate.

- **Interoperability:** The system should seamlessly integrate with various Electronic Health Record (EHR) systems and healthcare standards.

- **Maintainability:** The system must allow for easy updates, maintenance, and continuous improvement.

- **Compliance:** Strict adherence to healthcare data privacy regulations, such as HIPAA and GDPR, is essential.

- **Scalability:** The system should scale to handle increasing data volumes and user loads as healthcare institutions grow.

- **Performance:** The system should deliver fast OCR and NLP processing, ensuring efficient data management.

- **Customization:** Healthcare institutions should be able to customize the system to meet their unique needs.

## 4.4  Non-Functional Requirements

- **Performance:** The system shall process OCR and NLP tasks within a maximum response time of 2 seconds.

- **Availability:** The system should have a minimum uptime of 99.9 for healthcare professionals' continuous access.

- **Security:** Data transmitted between the system and users shall be encrypted using industry-standard encryption protocols.

- **Scalability:** The system should support a minimum of 100 concurrent users without performance degradation.

- **Usability:** The system's user interface shall comply with accessibility standards to accommodate users with disabilities.

- **Regulatory Compliance:** The system shall continuously update to remain compliant with evolving healthcare data privacy and security regulations.

## 4.5   Use Cases/ Use Case Diagram



Figure 4.1: Use Case Diagram

## 4.6    Sequence Diagrams/System Sequence Diagram



Figure 4.2: Sequence Diagram

## 4.7 Test Plan (Test Level, Testing Techniques)



Figure 4.3: Test Plan

## 4.8   Software Development Plan

**SEP-OCT**

**STAGE 1: STARTING**
- Initial Research+Literature Review
- Data Gathering
- Annotations
- Labelling
- Text Detection Techniques

**JAN-MAR**

**STAGE 3: TRAINING AND TESTING**
- Completing FRONT-END and BACK-END Web Application
- Improving Text Extraction Module
- Testing And Validation
- Integration Of Improved Detection Techniques
- User Testing And Feedback

**STAGE 2: BASIC IMPLEMENTATION**
- Detailed Literature Review
- Analysis
- Training Text Extraction Models
- Improving Text Detection Techniques
- Creating FrontEnd and BackEnd (Initial Stage)

**STAGE 4: FINALIZING**
- Documentation
- Training And Support
- Performance Optimization
- Deployment
- Project Presentation And Final Evaluation

**NOV-DEC**

**APR-MAY**

Figure 4.4: Software Development Plan

# 4.9 Architecture Diagrams



Figure 4.5: Architecture Design

# Chapter 5

# Iteration Plan

**SEP-OCT**

**STAGE 1: STARTING**

• Initial Research+Literature Review
• Data Gathering
• Annotations
• Labelling
• Text Detection Techniques

**JAN-MAR**

**STAGE 3: TRAINING AND TESTING**

• Completing FRONT-END and BACK-END Web Application
• Improving Text Extraction Module
• Testing And Validation
• Integration Of Improved Detection Techniques
• User Testing And Feedback

**STAGE 2: BASIC IMPLEMENTATION**

• Detailed Literature Review
• Analysis
• Training Text Extraction Models
• Improving Text Detection Techniques
• Creating FrontEnd and BackEnd (Initial Stage)

**STAGE 4: FINALIZING**

• Documentation
• Training And Support
• Performance Optimization
• Deployment
• Project Presentation And Final Evaluation

**NOV-DEC**

**APR-MAY**

# Chapter 6

# Designs

## 6.1 Architectural design



Figure 6.1: Architecture Design

# 6.2   Component Diagram



Figure 6.2: Component Design

## 6.3 Layer Diagram



Figure 6.3: Layer Design

## 6.4 Flow Diagram



Figure 6.4: Flow design

## 6.5 Activity Diagram



Figure 6.5: Activity Design

# Chapter 7

# Iteration 3

The iteration is expected to be completed by the midterm of the FYP-2. This chapter will have some of the artifacts based on system design. The requirements analysis section is same for all the systems while the design may vary. There may have two types of designs the structural design or . First section is for the structural design.

structural design

## 7.1 Component Diagram



Figure 7.1: Component Diagram

## 7.2  Layer Diagram

Below is the layer diagram representing each working layer of our system in a systematic way.



START

USER CHOOSES MODEL

USER PROVIDES DOCUMENT

CHECK FORMAT

INITIAL PROCESSING

APPLY ML MODELS

SEARCH TEXT IN MEDICINE DATABASE

FORMATTING PRESCRIPTION

DISPLAYING TO THE USER

STOP

## 7.3   Flow Diagram



Figure 7.3: Flow Diagram

## 7.4  Data Flow Diagram (DFD)



Figure 7.4: Data Flow Diagram

## 7.5  Data Dictionary

In this project we are using various available OCR datasets. We've scrapped multiple platforms for this purpose after it, we've fine tunned our model for the best performance.

# 7.6 Activity Diagram



Figure 7.5: Activity Design

## 7.7 Network Automata/ Graphs or State Machine



Figure 7.6: Network Automata

## 7.8   Call Graph or Sequence Diagram



Figure 7.7: Sequence Diagram

## 7.9   Use Case Diagram



Figure 7.8: Use Case Diagram

## 7.10    Algorithm Design



Figure 7.9: Algorithm Design

## 7.11    Development Phase



Figure 7.10: Development Phase

### 7.11.1 Suites or Test Cases

```
8]:  hw_image1 = hw_image.crop((0, 250, hw_image.size[0],370))
     display(hw_image1)
```

virtually indistinguishable from actual
handwriting.

```
9]:  ocr_image(hw_image1)
```

```
9]:  'virtually indistinguishable from actual'
```

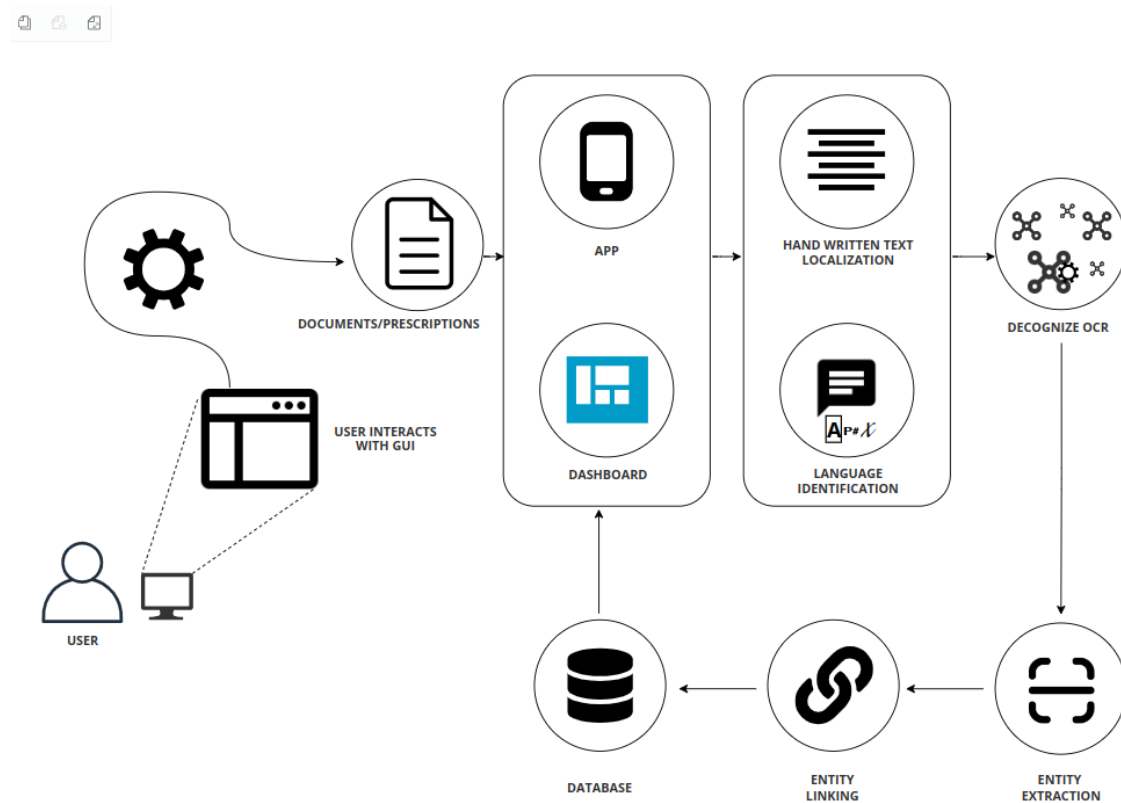Figure 7.11: Suites or Test Cases

### 7.11.2 Deployment Diagram



Figure 7.12: Deployment Diagram

## 7.12 Test Plan (Test Level, Testing Techniques)



Figure 7.13: Test Plan

### 7.12.1 SVN or GitHub (Optional)

You can find the code for Decognize on GitHub at:
https://github.com/mahadashraf/Decognize$_{ocr}$

### 7.12.2 Configuration/ Setup and Tool Manual (Optional)

Initially you have to install any python working platform. Then you've to install all the dependencies, including the libraries e.t.c. Then you've to integrate our trained model with the interface and knowledge graph. After this step your system in now fully functional. Input your document and accordingly test out the output.