

SUMMER INTERNSHIP REPORT

Presented at:

**National School of Electronics and Telecommunications of
Sfax**

Submitted by:
Imen Bakir

**Development of a Web application to detect fraud in the
pharmaceutical packaging**

Completed within:

Pixemantic



Supervised by:
Mr TOUATI Mohamed

Email:
mohamed.touati@pixemantic.com

Acknowledgements

I would like to extend my sincerest thanks to my supervisor, **Mr. Mohamed Touati**, for providing me with the opportunity to undertake this enriching internship at Pixemantic. This experience has been a true revelation, allowing me to explore my skills and deepen my interest in artificial intelligence and machine learning.

I also want to express my gratitude to the entire **Pixemantic** team. Their availability and willingness to provide detailed answers to each of my questions greatly contributed to my learning. Their kindness and patience have deeply impressed me, and I am immensely grateful to them.

I would also like to extend a warm thank you to my academic institution, **ENET'COM**, which not only provided me with the necessary knowledge but also opened the doors to professional experience. I also thank the jury for taking the time to evaluate my modest work.

Finally, I want to express my gratitude to all the people who, directly or indirectly, contributed to the completion of this work. Your support has been invaluable, and I am deeply grateful to you.

Contents

Acknowledgements	I
General Introduction	1
1 General Context	2
1.1 Introduction	2
1.2 Presentation of the Hosting Organization	2
1.2.1 Company Description	2
1.2.2 Sector of Activity	3
1.2.3 Introduction	3
1.2.4 Problem Statement	3
1.2.5 Proposed Solution	3
2 Technical Aspects	5
2.1 Introduction	5
2.2 Programming Languages	5
2.2.1 Python	5
2.2.2 HTML5 & CSS3	5
2.3 Frameworks/Libraries Used	6
2.3.1 BeautifulSoup 4	6
2.3.2 OpenCV	6
2.3.3 EasyOCR	6
2.3.4 SpaCy	7
2.3.5 SciKit-Learn	7
2.3.6 Pandas	7
2.3.7 Matplotlib	8
2.3.8 NLTK	8
2.3.9 Streamlit	8
2.4 Development Tools	9
2.4.1 Conclusion	10
3 Implementation	11
3.1 Introduction	11
3.2 Definitions	11
3.2.1 Optical Character Recognition (OCR)	11
3.2.2 Named Entity Recognition (NER)	11
3.3 Solution Architecture	11
3.4 Data Collection	12

Contents

3.5	Image Preprocessing	13
3.5.1	Enhancing Image Quality with ESRGAN	13
3.5.2	Orientation Correction	13
3.5.3	Text Extraction	14
3.5.4	Named Entity Recognition with SpaCy	14
3.5.5	Fraud Detection Using Jaccard Similarity Metric	15
3.6	Presentation of Our Solution	16
3.6.1	Authentication Interface	16
3.6.2	Registration Interface	16
3.6.3	Home Interface	16
3.6.4	Loading the Image and Extracting Text	18
3.6.5	Display of Labeled Text	19
3.6.6	Download PDF Summary	19
3.6.7	Fraud Detection	20
3.7	Conclusion	22
	General Conclusion	23

List of Figures

3.1	Solution Architecture.	12
3.2	Image Orientation Correction	13
3.3	medicale.json	14
3.4	Summary File	15
3.5	PharmaCare Logo	16
3.6	Authentication Interface	17
3.7	Registration Interface	17
3.8	Response to Membership Request	18
3.9	Home Interface	18
3.10	Load the Image from the Computer	19
3.11	Capture Image with Camera	19
3.12	Extracting Text	20
3.13	Display of Labeled Text	20
3.14	Summary File	21
3.15	Fraud Detection	21

List of acronyms and abbreviations

OCR *Optical Character Recognition*

NER *Named Entity Recognition*

HTML5 *HyperText Markup Language 5*

CSS3 *Cascading Style Sheet*

NLTK *Natural Language ToolKit*

General Introduction

In a world faced with an alarming crisis of counterfeit medicines in Sub-Saharan Africa, our unwavering commitment to combating pharmaceutical fraud has led to the creation of our innovative solution: "**PharmaCare**". The data provided by the **United Nations Office on Drugs and Crime** [1] is alarming, with nearly **half a million** estimated lives lost each year in the region due to counterfeit drugs. At the heart of this tragedy, **267,000** deaths are attributed to counterfeit or substandard antimalarial drugs, while **169,271** lives are at risk due to counterfeit or substandard antibiotics, especially in the context of severe childhood pneumonia.

In this concerning context, the 2023 report on the "**Traffic in Medical Products in the Sahel**" [2] paints a complex picture. Although the precise quantification of the extent of the traffic in medical products remains a challenge, studies converge to highlight the ongoing prevalence of counterfeit and substandard medical products, with estimates ranging from **19 to 50 percent**. International operations conducted between 2017 and 2021 seized no less than **605 tonnes** of various medical products in West Africa, underscoring the urgency of addressing this issue.

Sahel countries, including Burkina Faso, Chad, Mali, Mauritania, and Niger, largely rely on the importation of medical products due to a lack of developed local pharmaceutical infrastructure. The influx of counterfeit medicines primarily originates from pharmaceutical exporters, including Belgium, France, China, and India, which are diverted from legal supply chains or manufactured in neighboring nations.

Counterfeit medicines encompass those sold without approval, authorization, or a license and may be expired or lack essential active ingredients, thus exposing patients to significant risks when seeking medical care.

Our application, "**PharmaCare**," relies on the use of Optical Character Recognition (**OCR**) and Medical Named Entity Recognition (**NER**) techniques to detect fraud in drug packaging. By comparing the identified entities to our meticulously curated drug database, we empower users to verify the authenticity of their medicines. Our ultimate vision? A world where fraudulent medicines are eradicated, thereby preserving the health and integrity of every individual.

Chapter 1

General Context

1.1 Introduction

This chapter serves as an introduction to our internship project. This project requires a detailed study of certain concepts that not only pertain to the general framework of the project but also its implementation. The first part of this chapter will be dedicated to presenting the startup where I conducted my internship. As for the second part, it will focus on the overall idea of our project and the problem that drove us to develop this application.

1.2 Presentation of the Hosting Organization

1.2.1 Company Description



pixematic

Pixematic [4] is a startup specializing in the automatic interpretation of images and data from imaging sensors (laser, radar, infrared, etc.). It was founded in 2020, bringing together experts from Finest AI and consultants in risk management from Crédit Agricole Assurance, BNP, and IBM.

Pixematic was founded on our beliefs related to digitalization, AI, and data, as well as through ambitious projects that address business challenges in any industry or profession. We gather a community of recognized experts and passionate collaborators who work in an AGILE mode hand in hand with the client to deliver the project in a very short timeframe. We create value-generating solutions with our clients, improving their productivity, services, and thus their economic performance.

1.2.2 Sector of Activity

Pixementic was founded on convictions rooted in digitalization, AI, and data, as well as through ambitious projects addressing the business challenges of any company, regardless of its industry or profession. A community of recognized experts and passionate collaborators comes together at **Pixementic**, working in an AGILE mode alongside clients to ensure rapid project delivery. Value-generating solutions that improve clients' productivity, services, and therefore economic performance are developed in collaboration with them.

Project Presentation

1.2.3 Introduction

In a world where the global pharmaceutical industry is plagued by counterfeit drugs, our commitment to public safety has led to the creation of "**PharmaCare**," an innovative solution designed to protect consumers from the dangers of fraudulent medicines. According to alarming statistics from the World Health Organization [1], counterfeit drugs contribute to a significant loss of lives and pose a serious threat to public health. Each year, approximately half a million people lose their lives due to counterfeit drugs, with a substantial number of deaths linked to counterfeit or substandard drugs in Sub-Saharan Africa alone.

1.2.4 Problem Statement

Detecting fraud in the pharmaceutical industry from an image of a medicine package raises significant concerns, as counterfeit drugs represent a considerable danger to global public health. How can we effectively extract and analyze text from these images to identify signs of counterfeiting, thereby helping to minimize the risks associated with fraudulent medicines and ensure patient safety?

1.2.5 Proposed Solution

Our solution relies on cutting-edge Optical Character Recognition (**OCR**) and Named Entity Recognition (**NER**) technology. Initially, we extract text from medicine package images using **OCR**. Next, we automatically identify and label key information from the packaging, such as drug names, dosages, compositions, and more. We then compare this information with a database containing officially authorized medicines. This comparison allows us to quickly detect inconsistencies, discrepancies, or anomalies that may indicate potential pharmaceutical fraud. Our solution thus contributes to enhancing the safety and legitimacy of medicines for consumers.

Conclusion

Our application, "**PharmaCare**," based on **OCR** and **NER** technology, represents a crucial advancement for public health safety. By ensuring the authenticity of medicines, we contribute to minimizing the risks associated with fraudulent drugs, which is essential for the health and well-being of all.

Chapter 2

Technical Aspects

2.1 Introduction

After presenting our project **”PharmaCare”** aimed at detecting pharmaceutical frauds in the previous chapters, we now delve into the technical aspect of our solution. This chapter is dedicated to describing the essential tools and technologies that played a crucial role in the successful development and implementation of our application. We will highlight the fundamental technological components that contributed to the realization of our project.

2.2 Programming Languages

The programming languages used for the development of our application are as follows:

2.2.1 Python



Python [5] is a high-level interpreted general-purpose programming language. Its design philosophy emphasizes code readability through the use of indentation. Python is also characterized by dynamic typing and automatic memory management.

2.2.2 HTML5 & CSS3



HTML5 is not just the successor to HTML 4; it is much more than that. While HTML 4 and other XHTML focused solely on web page content, HTML5 focuses on web applications and interactivity, without neglecting accessibility and semantics.



CSS3's basic principle is to separate the content of the page from its appearance. The HTML page contains the information, not how the information is displayed. Multiple displays are possible for a single content, such as monochrome displays, small screens, oral rendering (the web page content is read aloud), printing on paper, printing on transparencies, and more.

2.3 Frameworks/Libraries Used

2.3.1 Beautiful Soup 4

BeautifulSoup⁴ **Beautiful Soup 4** is a powerful and flexible Python library used for extracting data from HTML and XML documents. It provides user-friendly methods and features for easily navigating and manipulating the content of these documents, making it an essential tool for web scraping and data analysis from web pages.

2.3.2 OpenCV



OpenCV is a library of functions primarily designed for real-time computer vision. It was initially developed by Intel, then taken over by Willow Garage and Itseez. This library is cross-platform and freely usable under the Apache 2 open-source license.

2.3.3 EasyOCR



EasyOCR is a Python library that enables computer vision developers to easily perform Optical Character Recognition (OCR) seamlessly.

2.3.4 SpaCy

spaCy

SpaCy is an open-source natural language processing library with a fast statistical entity recognition system. The NER methods available in SpaCy assign labels to textual data and classify them. SpaCy also offers the ability to add arbitrary classes to entity recognition systems and update the model to include new examples. We can train our own data for business-specific needs and prepare the model as needed.

2.3.5 SciKit-Learn



Scikit-Learn is an open-source machine learning library in Python. It powers many AI and data science applications and provides simple and efficient tools for data exploration and analysis, accessible to everyone and reusable in various contexts. It is built on NumPy, SciPy, and matplotlib.

2.3.6 Pandas



Pandas is a software library written for the Python programming language for data manipulation and analysis. It offers data structures and operations for manipulating numerical tables and time series. It is released under a three-clause BSD license.

2.3.7 Matplotlib



Matplotlib is a plotting library for the Python programming language and its numerical mathematics extension, NumPy. It offers an object-oriented API for embedding plots into applications using general-purpose GUI toolkits such as Tkinter, wxPython, Qt, or GTK.

2.3.8 NLTK



NLTK is a leading platform for building Python programs to work with human language data. It provides easy-to-use interfaces to over 50 corpora and lexical resources, such as WordNet, as well as a suite of text processing libraries for classification, tokenization, stemming, tagging, parsing, and semantic reasoning, along with an active discussion forum.

2.3.9 Streamlit



Streamlit is a Python library designed for rapidly creating interactive web applications. It significantly simplifies the development process by allowing developers to quickly turn their Python scripts into user-friendly web applications without worrying about underlying complexity. Streamlit makes data visualization, prototyping, and interactive user interfaces easy, all with a simple and intuitive syntax.

2.4 Development Tools



Colab [3] is a free Jupyter notebook environment that works entirely in the cloud. More importantly, it requires no installation, and the notebooks you create can be edited simultaneously by team members, similar to how you edit documents in Google Docs. Colab supports many popular machine learning libraries that can be easily loaded into your notebook.



XAMPP[8] is a suite of software that allows you to set up a local web server, FTP server, and email server. It is a free, open-source, cross-platform web server solution that stands for "X" (referring to one of the four operating systems it runs on: Windows, Linux, macOS, and Solaris), "Apache" (the web server software), "MySQL" (a popular relational database management system), "PHP" (a server-side scripting language), and "Perl" (a general-purpose programming language). It is designed to create a local web server environment on your computer, primarily for web development and testing purposes.



phpMyAdmin is a web-based interface for managing and administering MySQL databases. It provides users with the ability to perform operations such as creating, modifying, and deleting databases and tables, as well as executing SQL queries, all through a user-friendly interface accessible via a web browser.



TeXMaker[7] is a powerful all-in-one text editor. It is a free and open-source LaTeX editor with support for Unicode, spell checking, auto-completion, and code folding. It integrates a built-in PDF viewer with synctex support and continuous mode display. Texmaker is cross-platform, running on Linux, macOS, and Windows.

2.4.1 Conclusion

In this chapter, we have explored in detail the technical aspects of our development environment, which were crucial for the establishment of PharmaCare. These technological components will form the foundation upon which we build our application in the following chapter.

Chapter 3

Implementation

3.1 Introduction

After delving deep into the technical foundations of our development environment in the previous chapter, we now embark on an essential step: the implementation of our project, “**PharmaCare**”. This chapter aims to describe all the steps we have taken to bring our project to life, guiding you through the process of creating “**PharmaCare**”.

3.2 Definitions

3.2.1 Optical Character Recognition (OCR)

Optical Character Recognition (**OCR**) is also known as text recognition. OCR software extracts and reuses data from scanned documents, camera images, and PDF files containing only images. The OCR software identifies the letters in the image, turns them into words, and then into sentences, enabling access to the original content for editing. It also eliminates the need for manual data entry.

3.2.2 Named Entity Recognition (NER)

Named Entity Recognition (**NER**) is a common challenge in Natural Language Processing (NLP) aimed at identifying and classifying named entities. These entities are real-world objects identified by a name, such as locations, people, countries, or organizations (e.g., Microsoft or Asia). Using raw or annotated data, NER labels and classifies these entities, using linguistic and statistical approaches to develop recognition systems. NER models identify entities in the text and categorize them into appropriate classes.

3.3 Solution Architecture

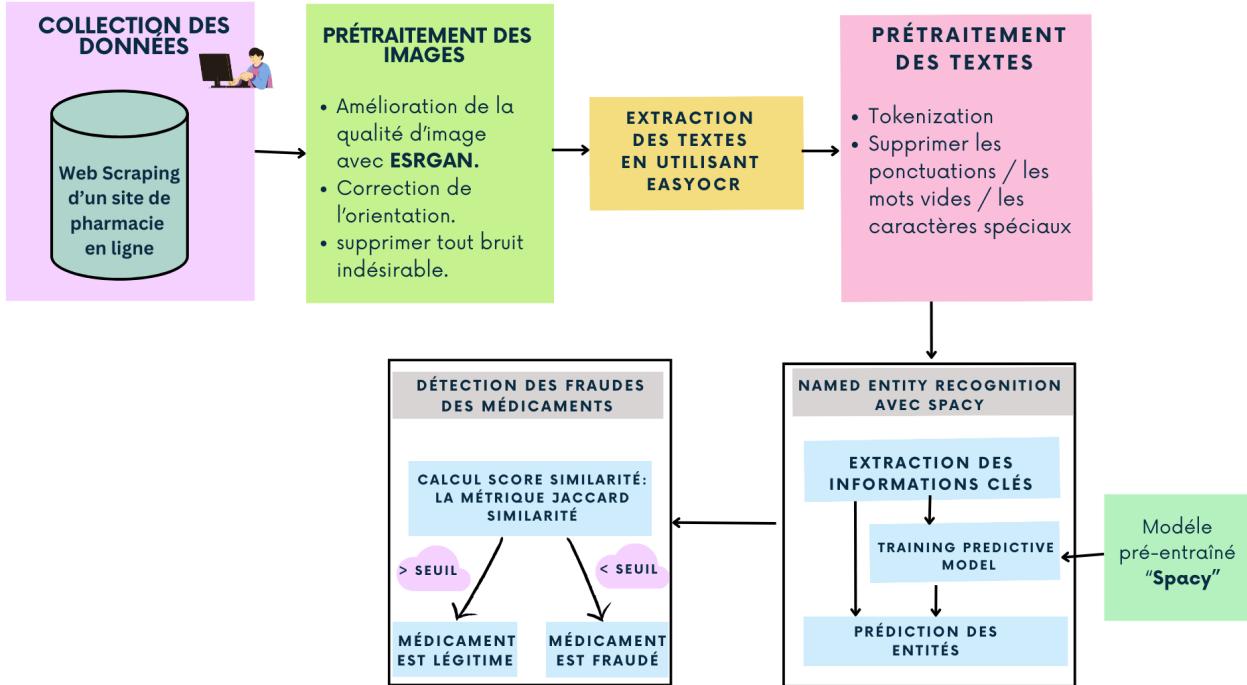


Figure 3.1: Solution Architecture.

Figure 3.1 illustrates the process of detecting pharmaceutical fraud in multiple stages. First, data is collected from an online medicine sales website[6]. Then, the images undergo rigorous preprocessing. This includes enhancing image quality through ESRGAN, correcting orientation, and removing any unwanted noise. Texts are then extracted from the images using EasyOCR. These extracted texts undergo additional preprocessing, including tokenization and the removal of punctuation, stopwords, and special characters. Named Entity Recognition (NER) with SpaCy is a crucial step, encompassing the extraction of key information and the training of a predictive model based on SpaCy. Finally, fraud detection relies on the Jaccard Similarity metric. If the score exceeds a specific threshold, the medicine is considered legitimate; otherwise, it is identified as fraudulent. We will describe each of these steps in detail to provide an in-depth understanding of how PharmaCare functions in detecting counterfeit medicines and ensuring consumer safety.

3.4 Data Collection

For fraud detection from a medical package image, the first crucial step is to collect data in the form of images. To accomplish this, we used the **Beautiful Soup** library to extract images from an online medicine sales website[6]. This data collection process forms the initial foundation of our fraud detection application, providing us with a diverse and comprehensive set of medical package images to analyze. Thanks to Beautiful Soup, we were able to efficiently navigate relevant web pages, extract images of medical packages, and build a rich and diverse database for our project. We also implemented the feature to capture images from the camera, allowing PharmaCare users to take pictures of medicines they want to verify.

3.5 Image Preprocessing

The second crucial step in implementing PharmaCare focuses on image preprocessing, an essential process to ensure the reliability of our application. This step was accomplished using a combination of advanced technologies, including the **ESRGAN** (Enhanced Super-Resolution Generative Adversarial Networks) model and the **OpenCV** library, to improve image quality, correct orientation, and remove any undesirable noise.

3.5.1 Enhancing Image Quality with ESRGAN

Our first task was to enhance the quality of images. We accomplished this mission by leveraging the **ESRGAN** model, an Enhanced Super-Resolution Generative Adversarial Network. **ESRGAN** was used to increase the resolution and sharpness of images, significantly improving the readability of text on medical packages. This ensures that the information subsequently extracted by our application is accurate and actionable.

3.5.2 Orientation Correction

Improper image orientation can pose significant challenges in text recognition. To overcome this challenge, we judiciously integrated the powerful **OpenCV** library into our process. Using **OpenCV**, we developed a sophisticated mechanism capable of automatically analyzing and correcting the orientation of each image. This step is crucial to ensure that the text is perfectly aligned and ready for optimal extraction. To illustrate the impact of this correction, we have included an example of an image before and after the orientation correction process.

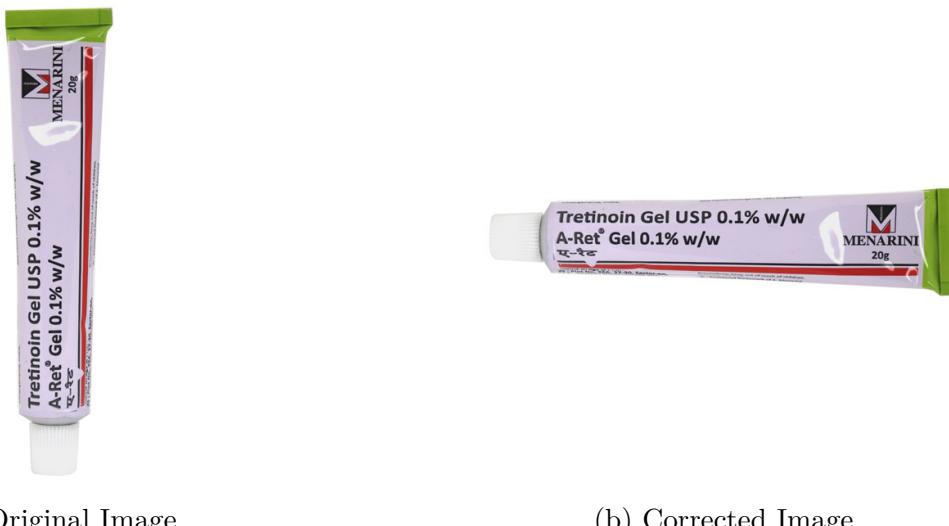


Figure 3.2: Image Orientation Correction

Thanks to this advanced preprocessing step, we have prepared our images for efficient and accurate Optical Character Recognition (OCR). The ESRGAN model, in collaboration with OpenCV, played a crucial role in improving the visual quality of the images,

allowing PharmaCare to extract textual information with high precision. This step forms the foundation on which our ability to successfully detect pharmaceutical fraud rests.

3.5.3 Text Extraction

Once the images are preprocessed, we used the **EasyOCR** library to extract text from each image. **EasyOCR** can automatically detect language and efficiently analyze text, which is crucial given the linguistic diversity of medical packaging. This text extraction step encompasses information such as the medicine's name, dosage, expiration date, batch number, etc., all of which are valuable for fraud detection.

3.5.4 Named Entity Recognition with SpaCy

This step was crucial in the development of PharmaCare. Thanks to this powerful natural language processing feature, we were able to accurately extract key information such as the medicine's name, type, composition, size, and dosage from the text extracted from medical packages. The ability to identify and classify these named entities is critical for our pharmaceutical fraud detection goal, as it allows us to specifically target relevant information. Furthermore, **SpaCy** provided us with the necessary flexibility to tailor the model to our specific domain by allowing us to add custom classes to better match the inherent variability in medical packaging.

The training of the **SpaCy NER** model relied on well-annotated medical entity data, including the "medicale.json" file. As shown in the figure below, the annotations included detailed information about medicines, types, compositions, sizes, and dosages, along with their exact positions in the text. This wealth of annotations allowed the model to effectively generalize from the provided examples, enhancing its ability to recognize similar entities in new medical packaging texts.



Figure 3.3: medicale.json

Once we extract named entities using SpaCy, PharmaCare generates a summary file containing these entities with their labels, as shown in Figure 3.4 below. This summary file simplifies the presentation of key information extracted from the medical package for the end user. It reinforces the user's confidence in the accuracy of the extracted data while facilitating verification of medicine details.

HEALTHCARE
A CLINICAL REPORT SUMMARY

COMPOSITION : Tretinoin Gel USP
DOSAGE : 0.1%
DRUGNAME : A-Ret
TYPE : Gel
SIZE : 20g

Figure 3.4: Summary File

In summary, thanks to Named Entity Recognition (NER) with SpaCy, PharmaCare can present crucial information in a clear and structured manner, such as composition, dosage, medicine name, type, and size, thereby enhancing consumer safety. This technological advancement enables PharmaCare to quickly detect pharmaceutical fraud, maintain high-quality standards, and inspire customer confidence.

3.5.5 Fraud Detection Using Jaccard Similarity Metric

The final step in our implementation process is to assess the reliability of the information extracted from medical images. To do this, we leverage the Jaccard Similarity metric, a set comparison method that allows us to measure the similarity between the information extracted from the loaded image and the labels of authentic medicines stored in the DataSet.

The Jaccard Similarity metric is calculated using the following formula:

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

Where $|A \cap B|$ represents the number of common entities between the two sets, and $|A \cup B|$ represents the total number of entities in those sets. The higher the similarity score, the more the information extracted from the image matches the authentic data from the reference dataset.

When the Jaccard similarity score is below a predefined threshold, it can be a potential indicator of fraud. This means that the information extracted from the image exhibits significant differences from the reference data, raising suspicions of counterfeiting or falsification.

This Jaccard Similarity-based approach significantly enhances PharmaCare's ability to detect pharmaceutical fraud accurately. It allows us to assess the consistency of the information extracted against legitimate reference data, thereby reducing the risks associated with potentially fraudulent medicines. Our application contributes to consumer safety by quickly identifying non-compliant products in the pharmaceutical market.

3.6 Presentation of Our Solution

In the following sections, we will provide an overview of the interfaces of our application, illustrating the various steps described above. Let's start by presenting the application's logo.



Figure 3.5: PharmaCare Logo

3.6.1 Authentication Interface

Each user has a private and secure area in our application. They must first enter their username and password to access their private space.

From this interface, if they are already registered, PharmaCare's user can log in. Simply enter their login and password and click on the "**Log In**" button to open their session. Of course, a new user interested in PharmaCare and wanting to take advantage of its pharmaceutical fraud detection features can sign up by carefully filling out the registration form. The figure below shows the new user registration interface on PharmaCare.

3.6.2 Registration Interface

After filling out the form, the internet user is informed that their request has been recorded. The Figure 3.8 below illustrates the response regarding the membership request.

3.6.3 Home Interface

After PharmaCare's user logs in, they are directed to the application's first interface, as illustrated in Figure 3.9. This interface contains a menu that allows quick access to various PharmaCare features:



Figure 3.6: Authentication Interface



Figure 3.7: Registration Interface

- **Extract text from image:** This option allows loading an image and extracting text.
- **Display labeled text:** It displays the extracted text with labels.

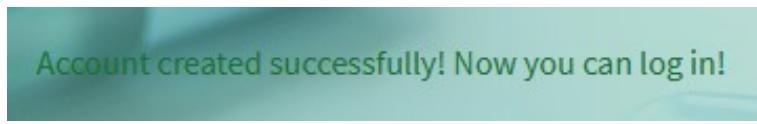


Figure 3.8: Response to Membership Request

- **Download PDF Summary:** For downloading a summary in PDF format.
- **Fraud detection:** For detecting pharmaceutical fraud.

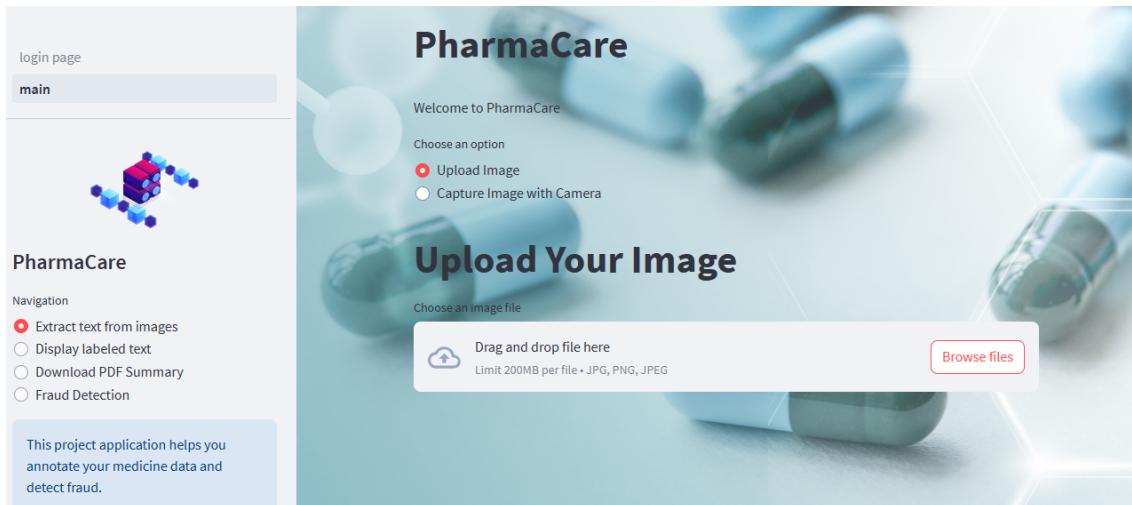


Figure 3.9: Home Interface

This menu simplifies the user's navigation within the application, allowing them to quickly access the desired functionality. By default, the application automatically directs the user to the **Extract text from image** function.

3.6.4 Loading the Image and Extracting Text

In this section, the user can choose to load an image of a medicine package to be tested using one of the following options:

1. Load the image from the computer:

By selecting the "**Upload Image**" option, the user can browse their computer files and choose the image to be tested (see Figure 3.10).

2. Capture the image directly from the computer's camera:

By choosing the "**Capture image with camera**" option, the user can activate the computer's camera and capture an image of the medicine package in real-time (see Figure 3.11).

Once the image is successfully loaded, the user can click on "**Extract text from image**" (see Figure 3.12) to initiate the process of extracting text from the image. This feature allows extracting essential information from medical packages for further analysis.

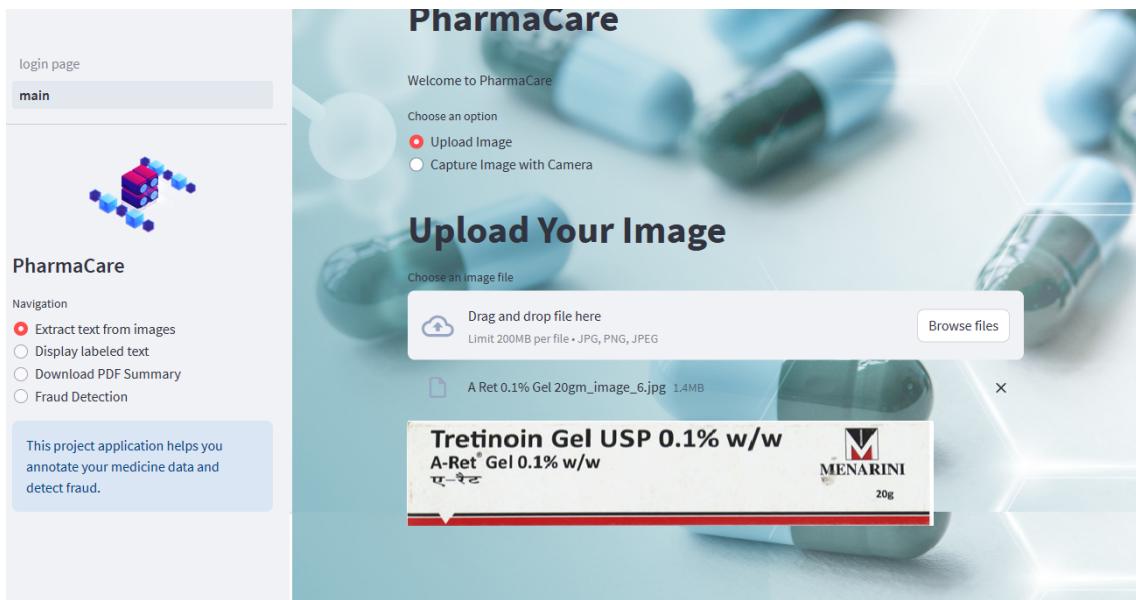


Figure 3.10: Load the Image from the Computer

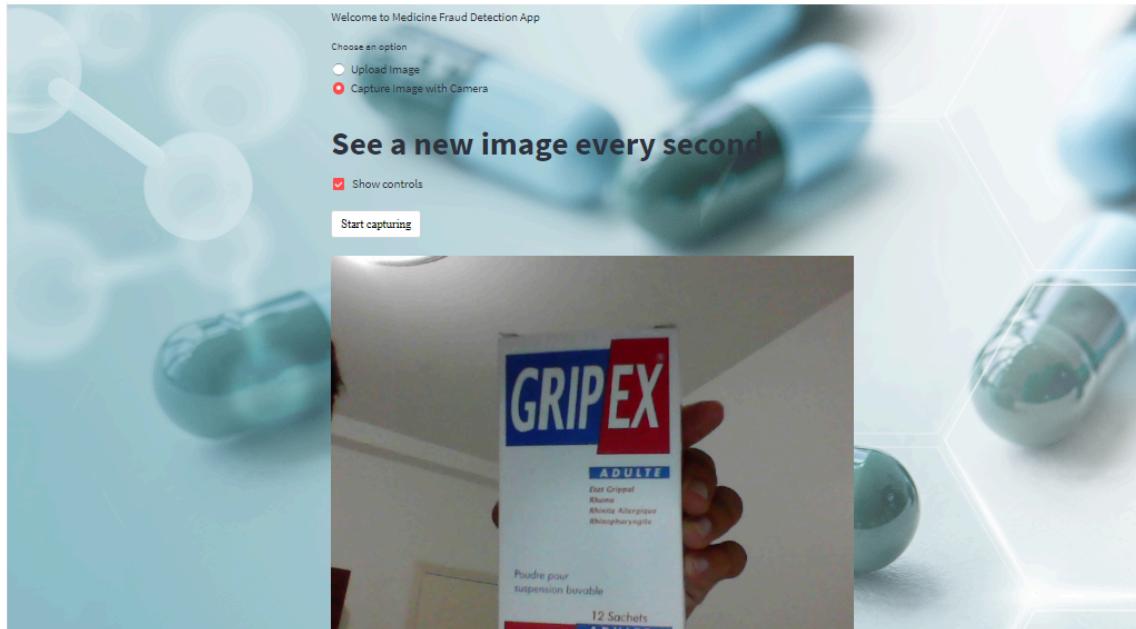


Figure 3.11: Capture Image with Camera

3.6.5 Display of Labeled Text

The user has the option to display the extracted text with the previously defined labels. To do this, they can select the **"Display labeled text"** option, as shown in Figure 3.13 below, where the labels will be visible.

3.6.6 Download PDF Summary

PharmaCare also allows users to download a concise summary in PDF format that includes the information extracted from the medicine image. Once we extract named entities

Chapter 3. Implementation

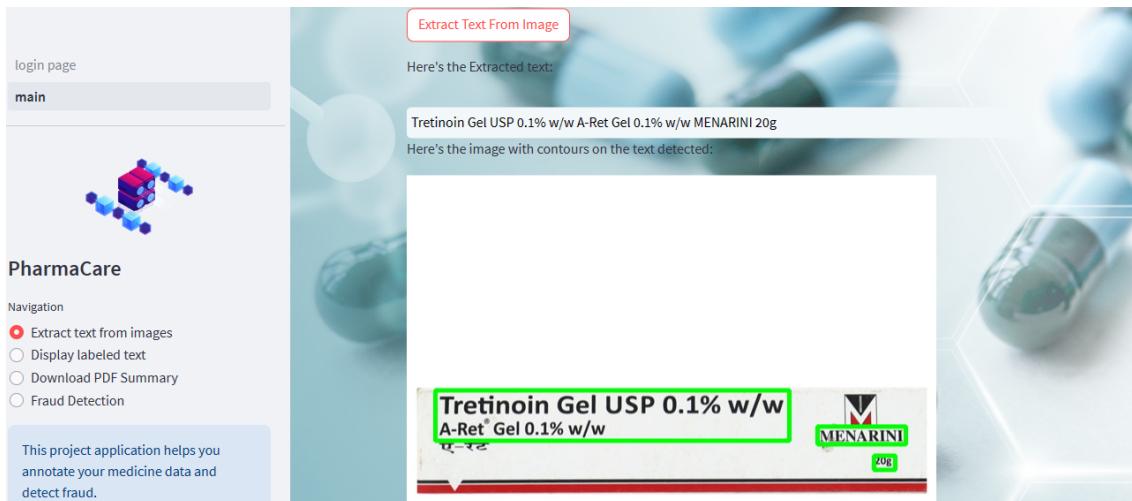


Figure 3.12: Extracting Text

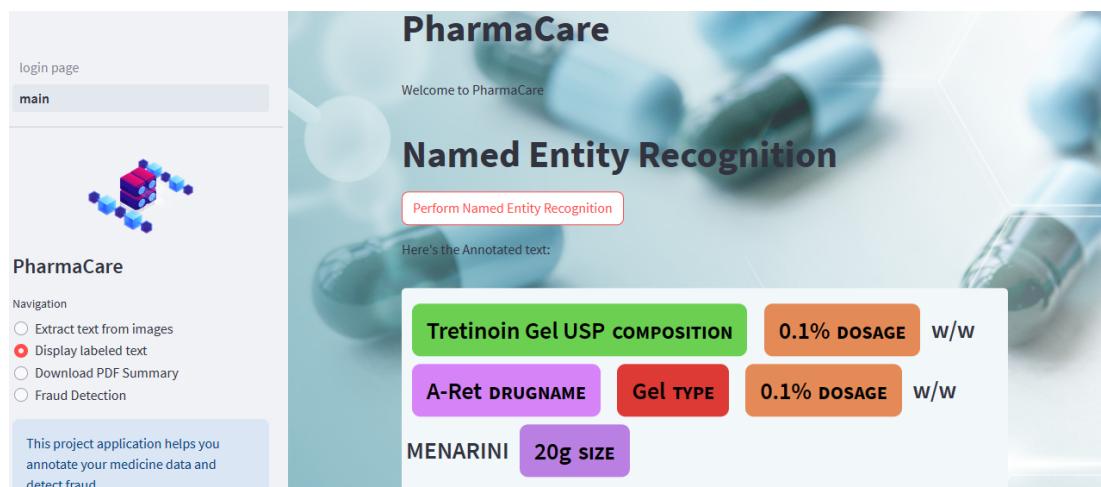


Figure 3.13: Display of Labeled Text

using SpaCy, PharmaCare generates a summary file containing these entities with their labels, as shown in Figure 3.14 below. This summary file simplifies the presentation of key information extracted from the medicine package for the end user. It reinforces the user's confidence in the accuracy of the extracted data while facilitating verification of medicine details. This step allows the user to explore the extracted information in a clear and organized manner. The **"Download PDF summary"** option (see Figure 3.14) facilitates this task by allowing the user to download this file for later use.

3.6.7 Fraud Detection

Finally, the main feature of our PharmaCare application is its ability to detect pharmaceutical fraud. When the user selects the **"Fraud detection"** option and clicks on the **"Calculation of max Jaccard Score"** button, they trigger the calculation of the similarity score between the information extracted from the loaded image and the labels of authentic medicines stored in the DataSet. This process also determines the fraud status associated with the medical package image (see Figure 3.15). This functionality is

Chapter 3. Implementation

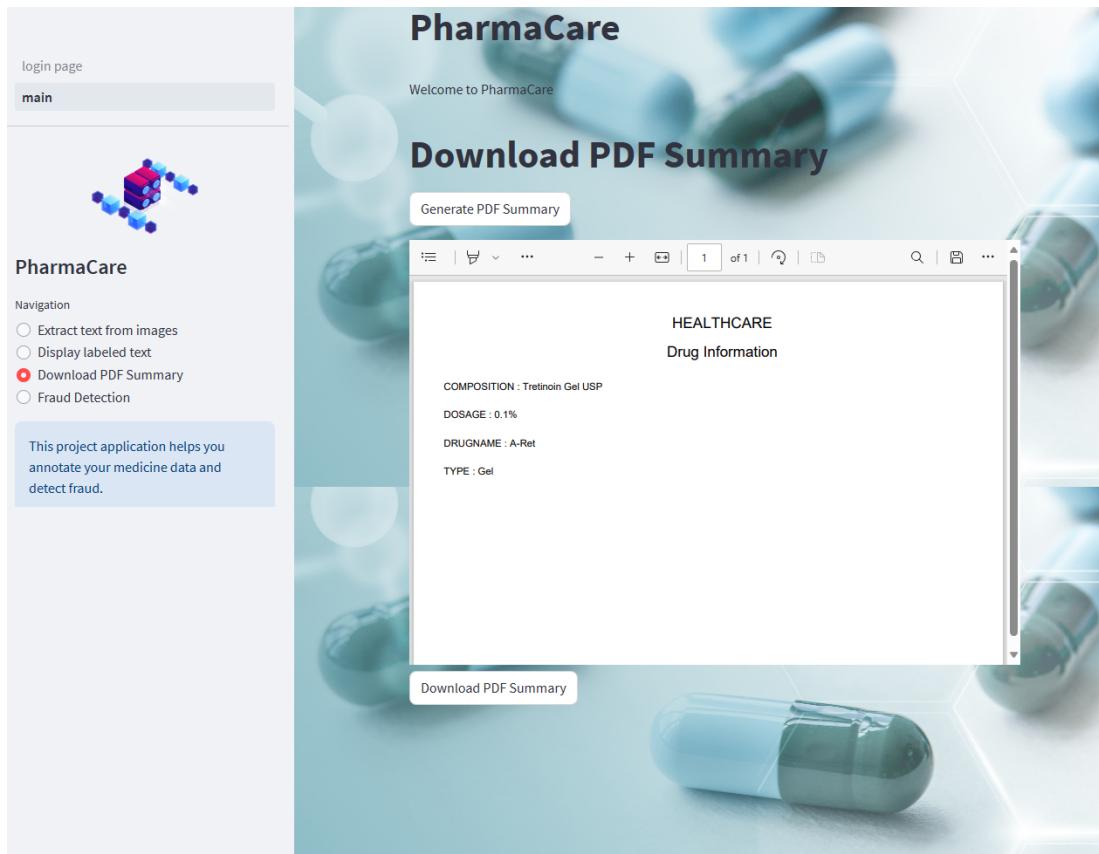


Figure 3.14: Summary File

at the core of our mission to ensure consumer safety by quickly and accurately identifying non-compliant pharmaceutical products on the market and thereby contributing to public health protection.



Figure 3.15: Fraud Detection

3.7 Conclusion

In conclusion, this introductory chapter allowed us to delve into the heart of the PharmaCare implementation process, our pharmaceutical fraud detection application. We followed an essential journey, starting with data collection from an online medicine sales website [6] and then moving on to advanced image preprocessing to ensure their quality and readability. We also explored the crucial steps of text extraction and named entity recognition (NER), which are at the core of PharmaCare’s ability to detect pharmaceutical fraud.

Finally, we presented PharmaCare’s user interface, which provides users with intuitive access to all these features. Whether it’s extracting text from an image, displaying labeled text, downloading a summary in PDF format, or performing fraud detection, our application offers a smooth and secure experience.

General Conclusion

The PharmaCare solution we have developed represents a crucial advancement in the fight against pharmaceutical fraud, particularly in sub-Saharan Africa, where counterfeit drugs tragically claim many lives each year. Our application relies on cutting-edge technologies, including Optical Character Recognition (OCR) and Named Entity Recognition (NER), to detect fraud in drug packaging and ensure the authenticity of pharmaceutical products.

The use of these technologies allows us to efficiently extract text from images of drug packages and automatically identify and label key information such as drug names, dosages, compositions, and more. This information is then compared to a database of legitimate drugs, enabling us to quickly detect inconsistencies, discrepancies, or anomalies that may indicate potential pharmaceutical fraud.

Our solution contributes to strengthening the safety and legitimacy of drugs for consumers, thereby reducing the risks associated with the consumption of fraudulent drugs. Using the Jaccard Similarity metric, we assess the reliability of extracted information and identify any significant deviations from reference data, which may indicate potential fraud. This approach ensures consumer safety by swiftly identifying non-compliant products in the pharmaceutical market.

In summary, PharmaCare is an innovative solution aimed at protecting public health by detecting and preventing pharmaceutical fraud. Our application serves as an essential tool in ensuring the authenticity of drugs, thereby helping save lives and preserving the health of every individual. It aligns with our commitment to combat a growing global issue and ensure consumer safety.

Webography

- [1] UNITED NATIONS OFFICE ON DRUGS and CRIME. *TRAFFICKING IN MEDICAL PRODUCTS IN THE SAHEL*. en. 2022. URL: https://www.unodc.org/documents/data-and%20analysis/tocta_sahel/T0CTA_Sahel_medical_2023.pdf.
- [2] Anna Fleck. *Up To 500,000 Killed by Fake Medicines in Sub-Saharan Africa*. en. May 24, 2023. URL: <https://cdn.statcdn.com/Infographic/images/normal/30068.jpeg>.
- [3] Google Colab. URL: <https://colab.research.google.com/>.
- [4] Pixementic. URL: <https://www.linkedin.com/company/pixemantic/?originalSubdomain=tn>.
- [5] Python. URL: <https://www.python.org/>.
- [6] Site de vente de médicaments en ligne. URL: <https://www.netmeds.com/prescriptions>.
- [7] TeXMaker. URL: <http://www.xm1math.net/texmaker/>.
- [8] XAMPP. URL: <https://www.apachefriends.org/download.html>.