

ENZIGMA- Assignment for AI Intern Position

ASSIGNMENT 01

Name- Akshay Bhosale

PRN- 202101040158

MIT Academy of engineering, Pune

Project Documentation: Automated Onboarding System with OCR and Database Integration

GitHub Link:- https://github.com/letsdoitbycode/OCR_detection

Demo drive links:-

https://drive.google.com/drive/folders/1JQt4TOC1ZH_4Xh9YR4IZYCMIM6eo1ORG?usp=sharing

Overview of the Solution

This project will develop an onboarding automation process, during which users can upload a scanned form in the format of an image. Solution: The textual information is captured from these forms, further consolidated into structured data, placed in a MySQL database, and provided with an easy-to-navigate user interface for uploading files as well as managing records.

Key Features:

1. File Upload: Upload scanned images in .jpg, .jpeg, or .png formats.
2. OCR Processing: Extract both printed and handwritten text from images.
3. Data Consolidation: Process and clean extracted text, converting it into a structured tabular format.
4. Database Storage: Save structured data into a MySQL database for easy management and retrieval.
5. Record Management:
 - View all records in a tabular format.

- Search records by name or email.

Steps to Set Up and Run the System

Prerequisites:

1. Python 3.x installed on your machine.
2. MySQL server running locally or on a remote server.

Environment Variables:

Create a .env file in the project root directory with the following variables:

```
1 //requirements.txt
2
3 groq
4 python-dotenv
5 numpy
6 streamlit
7 langchain-groq
8 langchain_community
9 pillow
10 mysql-connector-python
11 PyPDF2
12 PyMuPDF
13
```

```
1 GROQ_API_KEY = os.getenv("GROQ_API_KEY")
2 DB_HOST = os.getenv("DB_HOST")
3 DB_USER = os.getenv("DB_USER")
4 DB_PASSWORD = os.getenv("DB_PASSWORD")
5 DB_NAME = os.getenv("DB_NAME")
```

Running the Application:

1. Start the MySQL server and ensure the database is accessible.
2. Run the Streamlit application using the command:
➔ streamlit run app.py

Database Schema

Database Table: Record

The Record table stores extracted details in a structured format. Below is the schema:

Field	Type	Description
name	VARCHAR(255)	Full name of the candidate.

Field	Type	Description
email	VARCHAR(255)	Email address of the candidate.
address	TEXT	Residential address.
dob	DATE	Date of birth.
age	INT	Age of the candidate.
gender	VARCHAR(50)	Gender of the candidate.
mobile	VARCHAR(15)	Contact number.
education	TEXT	Educational qualifications.
profile	TEXT	Professional profile or job title.

Code and Functionalities

1. Image Upload and Preprocessing

The system allows users to upload scanned forms and preprocesses them for OCR. Preprocessing includes resizing the image and splitting it into horizontal stripes for better accuracy during text extraction.

```

1  # Function to encode an image to base64
2  def encode_image_pil(image: Image.Image) -> str:
3      buffered = io.BytesIO()
4      image = image.convert("RGB")
5      image.save(buffered, format="JPEG", quality=85)
6      return base64.b64encode(buffered.getvalue()).decode("utf-8")
7
8  # Function to split an image into horizontal stripes
9  def split_image_into_horizontal_stripes(image: Image.Image, stripe_count: int = 5, overlap: float = 0.1):
10     width, height = image.size
11     stripe_height = height // stripe_count
12     overlap_height = int(stripe_height * overlap)
13
14     stripes = []
15     for i in range(stripe_count):
16         upper = max(i * stripe_height - overlap_height, 0)
17         lower = min((i + 1) * stripe_height + overlap_height, height)
18         stripe = image.crop((0, upper, width, lower))
19         stripes.append(stripe)
20     return stripes

```

2. OCR Functionality

Extracts textual content from uploaded images. Both printed and handwritten elements are captured and processed.

```
1 # Function for OCR using Groq
2 def ocr(image: Image.Image, model: str = "llama-3.2-90b-vision-preview") -> str:
3     groq_llm = ChatGroq(
4         groq_api_key=GROQ_API_KEY,
5         model_name=model,
6         temperature=0
7     )
8
9     image_data_url = f"data:image/jpeg;base64,{encode_image_pil(image)}"
10
11     messages = [
12         {
13             "role": "user",
14             "content": [
15                 {"type": "text", "text": (
16                     "The uploaded image contains both printed text and handwritten notes. "
17                     "Your task is to carefully extract all textual content, including handwritten elements."
18                 )},
19                 {"type": "image_url", "image_url": {"url": image_data_url}}
20             ]
21         }
22     ]
23
24     response = groq_llm.invoke(messages)
25     return response.content.strip()
26
```

3. Data Consolidation

The extracted text from image sections is combined into a structured format, eliminating duplicates and resolving conflicts between overlapping sections.

```

1 # Function to consolidate markdown into tabular format
2 def format_to_table(markdown_runs: list, model: str = "llama-3.3-70b-versatile") -> str:
3     groq_llm = ChatGroq(
4         groq_api_key=GROQ_API_KEY,
5         model_name=model,
6         temperature=0
7     )
8
9     combined_markdown = "\n\n".join(markdown_runs)
10
11     messages = [
12         {
13             "role": "user",
14             "content": (
15                 "You are provided with multiple markdown outputs extracted from overlapping sections of an image."
16                 "Some sections may contain duplicate or conflicting information due to overlaps. "
17                 "Your task is to:"
18                 "\n\n1. Identify and consolidate rows of data that are related, ensuring that the most complete version of the information is retained."
19                 "\n2. For rows with conflicting information (e.g., different values for a field), prioritize the more detailed entry."
20                 "\n3. If a field is missing in one row but present in another, combine the information into a single row."
21                 "\n4. Output the consolidated data in a clean tabular format using Markdown syntax, suitable for direct rendering."
22                 "\n5. Output Only Markdown: Return solely the Markdown content without any additional explanations or comments."
23                 "\n\nHere is the data to process:\n\n"
24                 + combined_markdown
25             )
26         }
27     ]
28
29     response = groq_llm.invoke(messages)
30     return response.content.strip()
31

```

4. Data Parsing

Converts the consolidated tabular data into a Python dictionary format for database insertion.

```

1 # Function to parse the consolidated markdown into a dictionary for database insertion
2 def parse_markdown_to_dict(markdown_table: str) -> list:
3     rows = markdown_table.split("\n")[2:] # Skip the header row
4     records = []
5     for row in rows:
6         fields = [field.strip() for field in row.split("|")[1:-1]]
7         if len(fields) >= 9:
8             record = {
9                 "name": fields[0],
10                "email": fields[1],
11                "address": fields[2],
12                "dob": fields[3],
13                "age": int(fields[4]) if fields[4].isdigit() else None,
14                "gender": fields[5],
15                "mobile": fields[6],
16                "education": fields[7],
17                "profile": fields[8],
18            }
19            records.append(record)
20     return records

```

5. Database Operations

- Insert Records: Saves extracted details into the MySQL database.
- Fetch All Records: Retrieves all stored records.
- Search Records: Searches for records by name or email.

```

1  # Function to insert records into MySQL database
2  def insert_records_to_db(records: list):
3      try:
4          connection = mysql.connector.connect(
5              host=DB_HOST,
6              user=DB_USER,
7              password=DB_PASSWORD,
8              database=DB_NAME
9          )
10         cursor = connection.cursor()
11
12         query = """
13         INSERT INTO Record (name, email, address, dob, age, gender, mobile, education, profile)
14         VALUES (%s, %s, %s, %s, %s, %s, %s, %s, %s)
15         """
16         for record in records:
17             cursor.execute(query, (
18                 record["name"], record["email"], record["address"], record["dob"],
19                 record["age"], record["gender"], record["mobile"],
20                 record["education"], record["profile"]
21             ))
22         connection.commit()
23         st.success("Records have been successfully inserted into the database!")
24     except Error as e:
25         st.error(f"Error connecting to the database: {e}")
26     finally:
27         if connection.is_connected():
28             cursor.close()
29             connection.close()
30
31 # Function to fetch all records from the database
32 def fetch_all_records():
33     try:
34         connection = mysql.connector.connect(
35             host=DB_HOST,
36             user=DB_USER,
37             password=DB_PASSWORD,
38             database=DB_NAME
39         )
40         cursor = connection.cursor(dictionary=True)
41         cursor.execute("SELECT * FROM Record")
42         records = cursor.fetchall()
43         return records
44     except Error as e:
45         st.error(f"Error fetching data: {e}")
46         return []
47     finally:
48         if connection.is_connected():
49             cursor.close()
50             connection.close()
51
52 # Function to search records by name or email
53 def search_records(search_term):
54     try:
55         connection = mysql.connector.connect(
56             host=DB_HOST,
57             user=DB_USER,
58             password=DB_PASSWORD,
59             database=DB_NAME
60         )
61         cursor = connection.cursor(dictionary=True)
62         query = """
63         SELECT * FROM Record
64         WHERE name LIKE %s OR email LIKE %s
65         """
66         cursor.execute(query, (f"%{search_term}%", f"%{search_term}%"))
67         records = cursor.fetchall()
68         return records
69     except Error as e:
70         st.error(f"Error searching data: {e}")
71         return []
72     finally:
73         if connection.is_connected():
74             cursor.close()
75             connection.close()
76
77

```

6. Streamlit Interface

The user interface built with Streamlit includes:

1. Sidebar for Uploading Files: Allows users to upload images and displays a preview.
2. Main Section: Displays OCR results and buttons to save data to the database.
3. Record Management: Provides options to view all records or search specific entries.

Instructions for Using the User Interface

1. Uploading an Image:
 - Use the Upload Image section in the sidebar.
 - Preview the uploaded image to confirm its accuracy.
2. Processing and Saving Data:
 - After uploading, the image is processed to extract text.
 - Review the consolidated tabular data in the OCR and Results section.
 - Click Insert Records into Database to save the data.
3. Managing Records:
 - Use the View All Records button to see all stored entries.
 - Search for specific records using the Search Records input field.