MMC
Modeling
Computing
Mathematical

# Intelligent Automated System for Parsing and Ranking Resumes

Zahour O.[1], Sebbar A.[1,*], Zahour B.[2], Karim A.[1]

[1]*Information Technology and Modeling, Faculty of Sciences Ben M'sick,
Hassan II University, Casablanca, Morocco*
[2]*Faculty of Legal, Economic and Social Sciences, Ibn Zohr University of Agadir, Morocco*
*\*Corresponding author: abdosebbar205@gmail.com*

Resume parsing is a method used to extract key information from resumes, allowing for further actions such as candidate selection and ranking. In traditional recruitment processes, companies often handle thousands of resumes manually or require applicants to follow a pre-defined template. However, the evolving recruitment environment calls for more advanced technological solutions and efficient resume analysis methods. Although various basic techniques can analyze structured documents, they are inadequate for processing unstructured formats such as PDF, DOC, and DOCX. The current methods for resume parsing primarily rely on techniques such as BERT, Natural Language Processing (NLP), keyword-based models, and named entity recognition (NER) models. In response to this, the proposed system introduces a new approach that uses Computer Vision through YOLOv8 and Large Language Models (LLMs) for enhanced performance and broader API integration. YOLOv8 is used for resume segmentation, while Tesseract OCR extracts relevant information in variable text format. The extracted data are then processed by two LLMs using the Gemini and OpenAI APIs, which compute similarity scores and rank candidates according to specific criteria.

## 1. Introduction

In recent years, the growth of digital databases has made data processing increasingly challenging for organizations and companies. This created a need for tools that can efficiently manage large volumes of data. At the same time, recruiting suitable candidates is essential but remains a complex task for Human Resources (HR) departments. One of the key challenges is the selection and evaluation of candidates. Traditionally, HR departments rely on manual processes to analyze resumes, a time-consuming approach that is not only inefficient but also prone to human error.

This paper introduces an automated system designed for the Public Laboratory for Testing and Studies (LPEE) to address these challenges. The system streamlines resume processing, analysis, and classification by focusing on the skills of candidates, work experience, and educational background.

The proposed model utilizes Computer Vision techniques to analyze and filter resumes. Through models like YOLOv8, it divides certain sections of the resumes visually [1]. Tesseract OCR is used for text extraction from relevant image sections after segmentation [2]. Through methodology integration, it aids in the proper extraction of relevant details by recognizing and converting visual sections to formatted text for further processing.

Additionally, the system makes use of advanced methods that include integration of Large Language Models (LLM) through OpenAI and Gemini APIs to make the recruitment process automatic. LLMs that possess the capability to determine sophisticated context-based interconnectivities of words

https://orcid.org/0000-0002-4271-3303 (Zahour O.), https://orcid.org/0009-0000-6558-8560 (Sebbar A.),
https://orcid.org/0009-0006-7979-4430 (Zahour B.), https://orcid.org/0000-0002-6645-1141 (Karim A.)

provide richer and more specific text-based representations [3]. Using these models, similarities in job descriptions and resumes are measured by the system to rank applicants on their fit for jobs.

The objective of this paper is to present a solution to replace manual examination by extracting relevant categories from candidate resumes automatically. Our method identifies crucial work aspects such as experience, skills, and education from candidate resumes. Our extracted data are analyzed by our Large Language Models to identify how suitable candidate resumes are for offered jobs. Our models calculate similarity scores and rank applicants based on how suitable they are for the jobs offered to them.

## 2. Related work

The growing population of applicants has caused numerous applicants to apply for every offer of employment, so that recruiters have to sift through numerous applicants to identify suitable ones. Most research makes use of several of these detection and segmentation models that have been designed for these purposes. [4] employed YOLOv8 to identify helmet violations in real-time with very few annotations and with a high level of accuracy of 92.5%. This proves that YOLOv8 performs well in situations where small datasets have to be used. Similarly, in agricultural research, [5] demonstrated that through YOLOv8, plant leaves segmentation is efficiently performed with 89.3% accuracy in an autonomous plant growth monitoring environment. In aviation security research, [6] employed YOLOv8 to detect drones at an early stage at an impressive level of accuracy of 94.7%. All these research papers establish that more data makes performance higher in these complicated scenarios so that YOLOv8 performs exceptionally well in various complicated situations.

Optical Character Recognition (OCR) also plays an important role in extracting text from images. In the authors' work [7], the average detection error of Tesseract OCR was evaluated at 11.30%, with 153 words identified out of 173. However, the average error rate for identified words reached 67.65%, indicating room for improvement, particularly in English handwriting recognition. In contrast, in [8], recent improvements to Tesseract OCR for Tifinagh script recognition are highlighted, making it more robust in this context. Although Tesseract is effective in many cases, its performance varies significantly depending on the language and quality of the processed documents [7–9].

Large Language Models (LLMs) have revolutionized language analysis by enabling in-depth context awareness and coherent text generation [3]. LLMs possess the ability to analyze lengthy texts of varying levels of complexity, extract relevant details, and understand implicit nuances of words. LLMs in employment search processes analyze applicants' experience, skills, and education to map these to specifications of jobs to give detailed summaries [10]. They offer computing similarity scores and refine suggestions accordingly.

Recent studies [11, 12] have set LLMs relation extraction model performance at 85% in testing for evaluation. Besides that, [13] unveils that LLMs in named entity recognition have reduced classification errors by 15% to reach 92% in accuracy. All these results reveal how LLMs contribute significantly to streamlining text processing systems, such as the automatic generation of summaries.

## 3. Materials and methods

Finding the right personnel quickly and efficiently has become a major challenge for companies, especially when resources and time are limited. Identifying the most qualified candidates from a large number of resumes is increasingly time-consuming and requires significant staff resources due to the ever-growing population. To address this issue, we propose enhancing the overall preselection and selection process for the best candidates from a large pool of resumes. This will be achieved through the automation of the preselection and selection processes, allowing for a streamlined approach that ensures efficient analysis and processing of resumes, verifies candidate suitability for the position, and ultimately facilitates informed decision-making.

### 3.1. Data collection and preprocessing

The resumes used in this study were collected in various formats (PDF, DOC, and DOCX) by the Public Laboratory for Testing and Studies (LPEE). Our objective is to process and analyze these resumes, submitted by candidates applying for positions within the company, using innovative techniques based on Computer Vision and Large Language Models (LLMs). The dataset consists of 6 300 resumes representing diverse candidate profiles. We began by loading the dataset, which includes unstructured files in multiple formats (see Table 1), and converting all these papers to uniform format (images). Duplicate resumes were dropped while only retaining ones that contain all of the categories needed for us to analyze. Subsequently, the images were converted to grayscale to improve clarity by augmenting text-background contrast while minimizing noise caused by color inconsistency. Through preprocessing in this step, the model is in good shape to focus on text contours and shape more than before. After completing data preprocessing and cleansing, only 3 500 of the resumes remained to be used. They were split further into three subsets: 70% (approximately 2 450 YOLO-annotated resume images) for training, 20% (around 700 images) for validation, and the remaining 10% (around 350 images) for testing the model.

**Table 1.** File types comprising the candidate resumes dataset.

| File Type | Total |
|-----------|-------|
| PDF | 3959 |
| DOC | 1793 |
| DOCX | 548 |

### 3.2. Proposed method

The proposed solution first involves developing a Computer Vision-based system capable of analyzing and processing resumes from a large database of resumes in different formats and extracting the required information. Subsequently, the system integrates Large Language Models (LLMs) to provide advanced services, offering high-performance features such as similarity score calculation and classification. The approach of our solution is illustrated in Figure 1. The initial phase of the process involves the selection of resumes, which is part of the preparation phase. Companies receive many resumes in various formats, including PDF, DOC, and DOCX, for job offers.
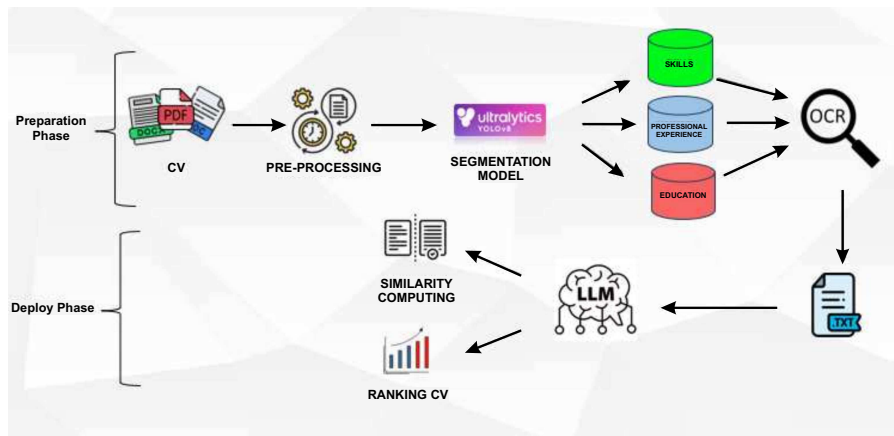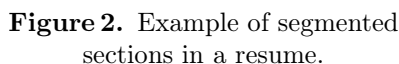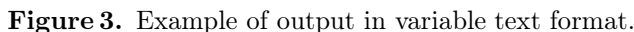


**Figure 1.** Architecture of the proposed model.

In our proposed model, we consider a set of resumes submitted in these formats as input. First, the input documents are transformed into a homogeneous format, specifically images. Next, we use the YOLOv8 model, which is based on a Computer Vision approach that detects and segments objects or specific sections from images, classifying them into predefined categories. In our case, these categories include skills, work experience, and education, which are extracted from resume images.

To streamline the category extraction process, we used Tesseract OCR to efficiently extract the text from the identified segments. The OCR process begins by dividing the image into parts corresponding to text blocks, followed by isolating lines, words, and finally characters. This allows us to identify the areas of the image containing text and structure the information into a usable format, making it easier to analyze while ensuring high accuracy, whatever the layout. Figure 2 illustrates how categories of interest are segmented in a resume image.

**Figure 2.** Example of segmented sections in a resume.

After extracting the entities, the information is stored in a structured text format (see Figure 3) in order to proceed to the deployment phase. At this stage, we employ Large Language Models (LLMs), based on high-performance architectures like the Transformer architecture, which allow us to effectively understand and manage long-term relationships between words in texts. These models are particularly efficient for tasks such as similarity score calculation, classification, and other advanced text analyses. One method used in our work is cosine similarity calculation, a powerful tool for measuring the similarity between two vectors in a vector representation space, formulated as follows [14]:

$$\cos(A, B) = \frac{A \cdot B}{\|A\|\|B\|}$$

This technique measures the angle between two vectors, making it possible to evaluate the semantic proximity between sentences or words.



**Figure 3.** Example of output in variable text format.

Recently, the accessibility of these Large Language Models via APIs provided by organizations has facilitated their integration into our processes. In this context, we chose to integrate two APIs: those of OpenAI and Gemini. This integration allows us to leverage their advanced capabilities to calculate the similarity score between a resume and a job offer, as well as to classify candidates based on their job designations, thereby optimizing our selection process. Subsequently, a performance comparison between the two APIs was conducted to identify the most reliable one in terms of similarity performance.

### 3.3. Performance metrics

In this study, several performance metrics were used to evaluate the effectiveness of our proposed model, namely: Precision, Recall, F1-score, and mean Average Precision (mAP). These metrics measure how well our model correctly detects relevant categories while minimizing errors. Each object class is evaluated independently, measuring the overlap area between predicted and reference zones.

— **Precision:** It measures the model's ability to correctly segment relevant sections of resumes (skills, education, and experience) while minimizing irrelevant inclusions. It is defined as the ratio of True Positives (TP) to the total number of True Positives and False Positives (FP) [15]:

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}.$$

— **Recall:** It evaluates the model's ability to identify all important resume sections. Instead of focusing on False Positives, it accounts for False Negatives (FN), representing cases where relevant information is missed [15]:

$$\text{Recall} = \frac{\text{TP}}{\text{TP} + \text{FN}}.$$

— **F1-score:** Since precision and recall often exhibit a trade-off, the F1-score is introduced as a harmonic mean of the two metrics. It provides a balanced evaluation of model performance [15]:

$$\text{F1-score} = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}.$$

— **Mean Average Precision (mAP):** To assess the model's ability across multiple classes, Mean Average Precision (mAP) is used. It represents the average of the average precision (AP) scores computed for each class, offering insight into overall segmentation performance [15]:

$$\text{mAP} = \frac{1}{N} \sum_{i=1}^{N} AP_i.$$

A high mAP indicates that the model is capable of accurately detecting relevant resume sections across various confidence levels, while a lower mAP suggests the need for further optimization.

## 4. Results and discussion

After applying our methodology, the obtained results (see Table 2) show the performance metrics of the customized YOLOv8 model (training model, skills model, and professional experience model) for detection and segmentation tasks. The metric measurements for the three models include an average precision of 95%, indicating that most positive predictions are correct, and an average recall of 94%, highlighting the models' ability to capture most of the present categories. The F1-score, which balances precision and recall, is 94%, confirming that the models are well-balanced and effective for detection. This demonstrates overall strong performance in terms of segmentation and accurate detection.

During the similarity score calculation and resume ranking phase for a given job offer, a comparative analysis of the similarity values obtained from the OpenAI and Gemini APIs was performed (see Figure 4). This step ensures a robust evaluation of candidate profiles based on specific criteria, offering reliable ranking results.

**Table 2.** Comparison of model performances across various metrics and their averages.

| Metric | Model 1 (Professional experience) | Model 2 (Skills) | Model 3 (Education) | Average |
|---|---|---|---|---|
| Precision | 0.9421 | 0.9503 | 0.9612 | 0.95 |
| Recall | 0.9305 | 0.9450 | 0.9600 | 0.94 |
| mAP50 | 0.9502 | 0.9523 | 0.9678 | 0.95 |
| mAP50-95 | 0.9020 | 0.9105 | 0.9204 | 0.91 |
| F1-score | 0.9363 | 0.9476 | 0.9606 | 0.9450 |

The results in Figure 4 present a comparison of similarity scores between OpenAI and Gemini for a set of 50 resumes. While both models show variations in performance, OpenAI consistently achieves higher similarity scores across most resumes. This consistency highlights its efficiency and reliability in identifying relevant matches for the job offer. In contrast, Gemini's performance is less uniform, with fluctuations suggesting varying levels of accuracy across different resumes.

Our findings align with previous studies [16], which observed similar trends in their analysis of both models, further confirming that OpenAI remains the most reliable choice for this task (see Figure 5).
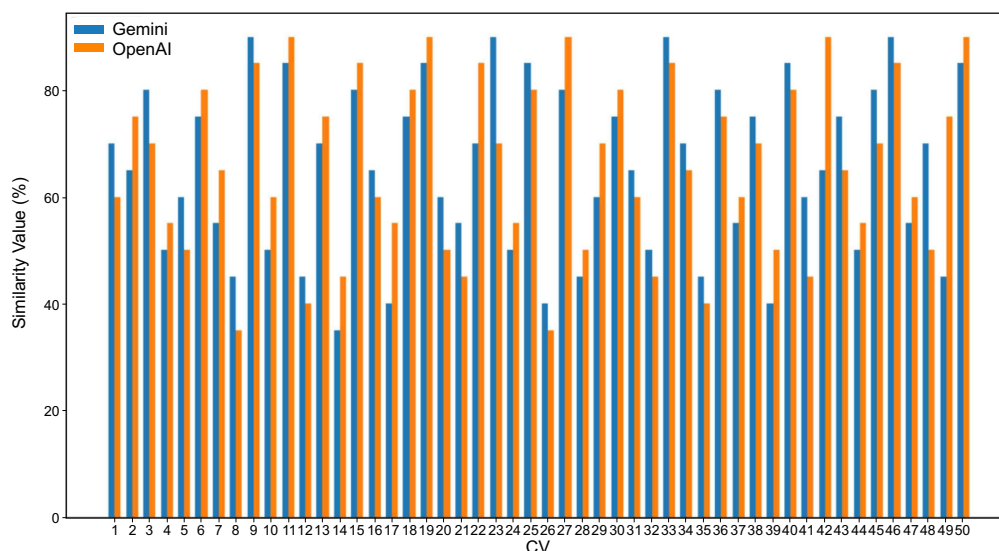
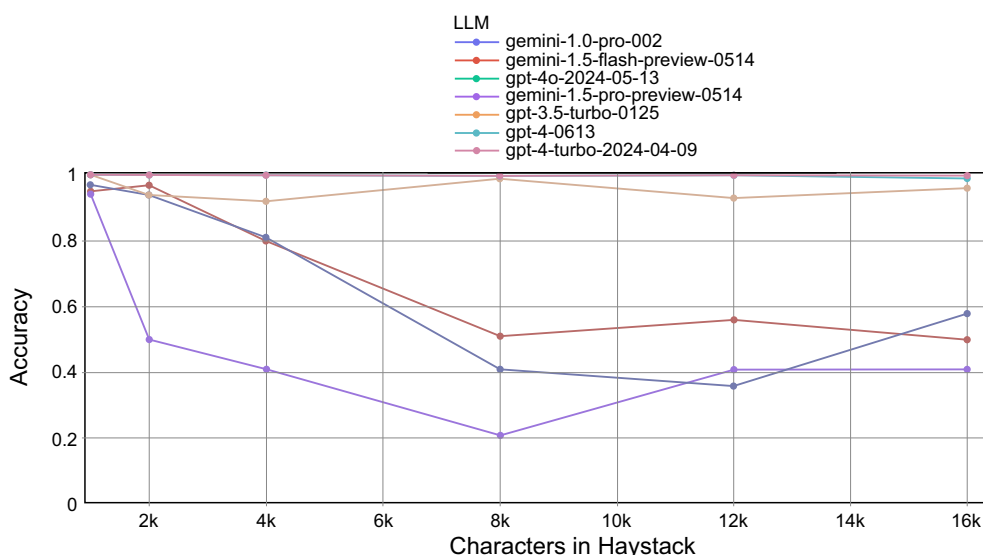**Figure 4.** Evaluation of similarity rates between Gemini and OpenAI.



**Figure 5.** Comparative assessment of LLM performances between OpenAI and Gemini [16].

## 5. Conclusion

In our work on resume processing, the model significantly improved the efficiency of the recruitment process by streamlining the analysis and ranking of candidates. Through the combined use of YOLOv8 for segmentation and Tesseract OCR for information extraction, manual tasks are reduced, thereby minimizing human error. The LLM models based on the Gemini and OpenAI APIs were highly effective in evaluating the match between resumes and job offers, with OpenAI showing slightly better performance in terms of consistency and overall efficiency. This automated system represents a major advancement for companies, providing a fast and accurate solution for managing applications, and it also has the potential to support career guidance for job seekers [17].

However, it is worth noting that our automated system has certain limitations. For example, access to the ChatGPT API is subscription-based, which may pose a financial constraint for some companies. Additionally, the system is limited by the API's request capacity, allowing only a specific number of resumes to be processed. This limitation may affect model scalability, particularly for large-scale recruitment operations that need to handle high volumes of resumes efficiently. Given these limitations, the Meta LLaMA model could present a promising research perspective in this field.

[1] Kang J., Zhao L., Wang K., Zhang K. Research on an improved YOLOv8 image segmentation model for crop pests. Advances in Computer, Signals and Systems. **7** (3), 1–8 (2023).

[2] Garai S. K., Paul O., Dey U., Ghoshal S., Biswas N., Mondal D. S. A Novel Method for Image to Text Extraction Using Tesseract-OCR. American Journal of Electronics & Communication. **3** (2), 8–11 (2024).

[3] Silva K., Frommholz I., Can B., Blain F., Sarwar R., Ugolini L. Forged-GAN-BERT: Authorship Attribution for LLM-Generated Forged Novels. Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics: Student Research Workshop. 325–337 (2024).

[4] Aboah A., Wang B., Bagci U., Adu-Gyamfi Y. Real-Time Multi-Class Helmet Violation Detection Using Few-Shot Data Sampling Technique and YOLOv8. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 5350–5358 (2023).

[5] Wang P., Deng H., Guo J., Ji S., Meng D., Bao J., Zuo P. Leaf Segmentation Using Modified YOLOv8-Seg Models. Life. **14** (6), 780 (2024).

[6] Yilmaz B., Kutbay U. YOLOv8 Based Drone Detection: Performance Analysis and Optimization. Computers. **13** (9), 234 (2024).

[7] Joshi K. Study of Tesseract OCR. GLS KALP: Journal of Multidisciplinary Studies. **1** (2), 41–50 (2021).

[8] Benaissa A., Bahri A., Allaoui A. E., Salahddine M. A. Build a Trained Data of Tesseract OCR Engine for Tifinagh Script Recognition. Data and Metadata. **2**, 185–185 (2023).

[9] Sporici D., Cuşnir E., Boiangiu C.-A. Improving the Accuracy of Tesseract 4.0 OCR Engine Using Convolution-Based Preprocessing. Symmetry. **12** (5), 715 (2020).

[10] Xu S., Wu Z., Zhao H., Shu P., Liu Z., Liao W., Li S., Sikora A., Liu T., Li X. Reasoning before Comparison: LLM-Enhanced Semantic Similarity Metrics for Domain Specialized Text Analysis. Preprint arXiv:2402.11398 (2024).

[11] Yahiaoui F., Limouni E. Génération et annotation de corpus pour l'entrainement et l'évaluation de modèles d'extraction de relations: utilisation de bibliothèques de génération de données et de LLMs. Institut des sciences informatiques et de leurs interactions – CNRS Sciences informatiques. hal-04678383 (2024).

[12] Résumé automatique de textes d'enquêtes judiciaires: retour d'expérience. Institut des sciences informatiques et de leurs interactions – CNRS Sciences informatiques. hal-04678366 (2024).

[13] De Murcia G., El-Allali I., Meineri L., Gillard L., Lastmann S. Rapport de Participation de Smart Tribune à EvalLLM2024: Quelques Usages de LLMs dans l'Univers de la Reconnaissance d'Entités Nommées. Atelier sur l'evaluation des modeles generatifs (LLM) et challenge d'extraction d'information few-shot. hal-04678371 (2024).

[14] Octavany O., Wicaksana A. Cleveree: an artificially intelligent web service for Jacob voice chatbot. TELKOMNIKA Telecommunication, Computing, Electronics and Control. **18** (3), 1422–1432 (2020).

[15] Anwar A. What is Average Precision in Object Detection & Localization Algorithms and how to calculate it? `https://is.gd/Jbgn1P`.

[16] Wiik L. OpenAI's GPT-4o vs. Gemini 1.5 Context Memory Evaluation. Medium (2024).

[17] Zahour O., Benlahmar E. H., Eddaoui A., Ouchra H., Hourrane O. A system for educational and vocational guidance in Morocco: Chatbot E-Orientation. Procedia Computer Science. **175**, 554–559 (2020).

# Інтелектуальна автоматизована система для розбору та ранжування резюме

Захур О.[1], Себбар А.[1], Захур Б.[2], Карім А.[1]

[1] *Інформаційні технології та моделювання, факультет природничих наук Бен Мсік,*
*Університет Хасана II, Касабланка, Марокко*
[2] *Факультет юридичних, економічних та соціальних наук,*
*Університет Ібн Зора в Агадірі, Марокко*

Парсинг резюме – це метод, який використовується для вилучення ключової інформації з резюме, що дозволяє здійснювати подальші дії, такі як відбір кандидатів та ранжування. У традиційних процесах рекрутингу компанії часто обробляють тисячі резюме вручну або вимагають від кандидатів дотримуватися попередньо визначеного шаблону. Однак, середовище рекрутингу, що постійно розвивається, вимагає більш просунутих технологічних рішень та ефективних методів аналізу резюме. Хоча різні базові методи можуть аналізувати структуровані документи, вони неадекватні для обробки неструктурованих форматів, таких як PDF, DOC та DOCX. Поточні методи парсингу резюме в основному спираються на такі методи, як BERT, обробка природної мови (NLP), моделі на основі ключових слів та моделі розпізнавання іменованих сутностей (NER). У відповідь на це запропонована система впроваджує новий підхід, який використовує комп'ютерний зір через YOLOv8 та моделі великих мов (LLM) для підвищення продуктивності та ширшої інтеграції API. YOLOv8 використовується для сегментації резюме, тоді як Tesseract OCR витягує відповідну інформацію у змінному текстовому форматі. Потім витягнуті дані обробляються двома LLM за допомогою API Gemini та OpenAI, які обчислюють оцінки подібності та ранжують кандидатів за певними критеріями.

**Ключові слова:** *резюме; комп'ютерний зір; YOLOv8; сегментація об'єктів; Tesseract OCR; LLM; оцінка подібності; класифікація; API; OpenAI.*