

# *Creating Text from Images with OCR API*

Mithila Prabhu  
 mithila.prabhu@stud.fra-uas.de

Khushal Singh  
 khushal.singh@stud.fra-uas.de

Taibaz Pathan  
 taibaz.pathan@stud.fra-uas.de

**Abstract**—Effective image preprocessing plays a crucial role in enhancing the accuracy and efficiency of Optical Character Recognition (OCR) systems by improving text clarity and reducing background noise. This paper presents a systematic evaluation of five essential preprocessing techniques grayscale conversion, adaptive thresholding, global thresholding, saturation adjustment and deskewing implemented using the Tesseract SDK. To facilitate this study, a C# console application is developed to apply these preprocessing methods and assess their impact on OCR performance. The evaluation is conducted based on three key factors: cosine similarity, time taken by each preprocessing technique and memory usage of each preprocessing technique. Cosine similarity metrics are employed to compare the extracted text against each preprocessing model, providing a quantitative measure of recognition accuracy. Additionally, time tracking is used to analyze processing speed, while memory usage is monitored to assess the computational overhead introduced by each technique. By systematically analyzing the trade-offs between accuracy and performance, this study offers valuable insights into optimizing preprocessing pipelines for real-world OCR applications.

**Keywords**—Optical Character Recognition (OCR), image preprocessing, Tesseract SDK, grayscale conversion, adaptive thresholding, global thresholding, saturation adjustment, deskewing, cosine similarity, processing time, memory usage.

## I. INTRODUCTION

Optical Character Recognition (OCR) is a foundational technology in the field of document analysis and recognition. It enables the automated extraction of textual information from scanned documents, PDFs, and natural images, thereby transforming unstructured visual content into machine-readable text. OCR has found widespread application in sectors such as banking, healthcare, education, and legal documentation, where it supports automation and improves operational efficiency. Traditional OCR engines, such as the Tesseract engine developed by Google, have relied on feature-based methods and heuristics for character recognition. Tesseract follows a two-stage pipeline: it first identifies connected components within an image to segment potential character regions and then applies a machine learning-based classifier to recognize the segmented text [1]. The integration of Long Short-Term Memory (LSTM) networks in later versions of Tesseract has substantially improved its performance, particularly in handling both printed and handwritten characters across diverse scripts. Despite these advancements, OCR systems continue to face challenges when processing noisy or low-quality images. Common issues include blurred text, varying font styles and

sizes, poor contrast, and complex or degraded backgrounds. To mitigate these effects, image preprocessing techniques are often employed as a preliminary step. Techniques such as binarization, noise filtering, edge detection, and adaptive thresholding have proven effective in enhancing the visibility and clarity of text regions prior to recognition [2]. In addition to improving image quality, recent research has explored the use of language models to enhance OCR output through contextual understanding. For example, Moon et al. proposed a hybrid approach that integrates BERT—a pre-trained transformer-based language model—with an OCR engine to refine noisy outputs by predicting masked words based on surrounding context [3]. This method employs cosine similarity to evaluate the visual similarity between the original and predicted word images, enabling more accurate word reconstruction under adverse imaging conditions. Beyond BERT, more advanced embedding techniques have emerged. OpenAI’s text-embedding-ada-002 model provides a lightweight yet powerful solution for generating semantic vector representations of textual content [4]. These embeddings can be employed in OCR post-processing pipelines for tasks such as duplicate detection, error correction, and semantic similarity evaluation. By incorporating such embeddings, OCR systems can better align predicted text with expected content, especially in domain-specific or context-sensitive applications. This research aims to develop an enhanced OCR pipeline by combining robust image preprocessing techniques with Tesseract OCR and embedding-based similarity scoring. The proposed system focuses on three key objectives: (1) improving text visibility through adaptive preprocessing; (2) evaluating multiple preprocessing pipelines to identify the most effective configuration; and (3) integrating semantic similarity measures using text embeddings to refine and validate OCR output. The system is evaluated on a diverse dataset containing images with varying levels of noise and distortion to demonstrate its practical utility and accuracy improvements over conventional approaches.

## II. LITERATURE REVIEW

Optical Character Recognition (OCR) systems in their early development primarily relied on character template matching techniques, which involved comparing an input image with a set of stored templates to identify the closest match. This approach proved effective for recognizing characters in printed text with consistent font types and sizes, particularly under

controlled conditions such as fixed grayscale bitmap images in known fonts like Times New Roman. However, the effectiveness of template matching significantly declined in scenarios involving variations in font style, image distortions, or poor image quality. The inability to generalize across different writing styles or adapt to noise and irregularities in the input made these early OCR systems less reliable in practical, real-world applications where such variability is common [5].

Recent advancements in deep learning, especially Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), have considerably improved the accuracy of OCR systems. Studies have shown that integrating CNN-based feature extraction with RNNs—such as Long Short-Term Memory (LSTM) networks—enhances OCR capabilities by capturing contextual dependencies between characters [6]. This combination has demonstrated substantial improvements in recognizing complex text structures and increasing accuracy in a variety of contexts.

Pre-processing is a fundamental step in improving the performance of OCR systems, particularly when dealing with low-quality inputs such as noisy, skewed, or poorly illuminated document images. Traditional pre-processing techniques—including grayscale conversion, global or adaptive thresholding, and contrast enhancement—have long been employed to increase the legibility of text prior to recognition. These methods serve to emphasize foreground text while suppressing irrelevant background information. In recent years, researchers have introduced more sophisticated pre-processing approaches that address specific image imperfections in a targeted manner. One such approach is proposed by Harraj and Raissouni, who designed a pre-processing pipeline tailored to scanned documents affected by uneven illumination. Their method integrates illumination correction with optimized grayscale transformation, followed by unsharp masking and global binarization. This combination not only enhances the visibility of characters but also improves the consistency of OCR outputs across varying lighting conditions. The authors demonstrated that this pre-processing strategy significantly boosts OCR accuracy, particularly for documents captured under suboptimal conditions [7].

Another critical aspect of OCR research is computational efficiency, especially in scenarios involving large-scale document processing. With the increasing deployment of OCR systems in real-world applications, optimizing both accuracy and processing speed has become essential. Singh, Colom, and Bontcheva conducted a comprehensive evaluation of OCR models on diverse datasets—including complex domains such as memes and document-style images—and emphasized the importance of balancing performance with computational demands. Their findings highlight that lightweight preprocessing strategies—particularly those involving grayscale conversion and adaptive thresholding—can enhance efficiency without significantly compromising accuracy, thereby making such approaches suitable for high-throughput OCR systems [8].

Misclassification of visually similar characters, such as ‘0’ and ‘O’ or ‘1’ and ‘l’, remains a persistent challenge in

OCR, particularly in off-line handwriting recognition. Arica and Yarman-Vural emphasize the importance of robust feature extraction and classification techniques to minimize such recognition errors. Their comprehensive review outlines how statistical, structural, and hybrid methods can be leveraged to improve accuracy by effectively distinguishing ambiguous characters during the classification phase [9].

The adoption of transformer-based models has brought notable advancements to OCR systems in recent years. One prominent example is the TrOCR model, which utilizes a combination of vision and language transformers pre-trained on large-scale datasets to recognize printed and handwritten text with high accuracy. This model eliminates the need for hand-crafted features by learning contextual and visual representations end-to-end. Empirical results have shown that TrOCR achieves superior performance on benchmark datasets, making it a strong candidate for tasks such as scene text recognition and digitization of historical documents [10].

Given these developments, the research presented in this paper builds upon prior work by integrating advanced image pre-processing methods, the Tesseract OCR engine, cosine similarity for text validation, and embedding-based text analysis, while also evaluating computational efficiency. The primary goal is to create a robust OCR pipeline that enhances text extraction accuracy while minimizing both time and memory consumption.

### III. METHODOLOGY

This section describes the methodology used in preprocessing images, extracting text, and evaluating the performance of different image processing models for OCR. The system consists of three major components:

- **Image Preprocessing** – Various transformations such as grayscale conversion, thresholding, saturation adjustment, deskewing, and shifting were applied to improve OCR accuracy.
- **OCR Processing** – The Terrasect SDK was used to extract text from preprocessed images.
- **Performance Evaluation** – The extracted text was analyzed based on cosine similarity, processing time, and memory usage to determine the best preprocessing model.

#### A. Image Preprocessing Techniques

Image preprocessing is crucial for enhancing text clarity before OCR. This study implements multiple transformations, each designed to address specific challenges such as noise, low contrast, skewed text, and misalignment.

1) *Grayscale Conversion*: Grayscale conversion is a crucial preprocessing step in Optical Character Recognition (OCR) that simplifies image complexity while preserving text integrity. By transforming a color image into a single-channel representation, grayscale conversion enhances contrast and prepares the image for further processing techniques, such as thresholding, which improves text clarity.

In this project, grayscale conversion is implemented using the `ConvertToGrayscale` class, which processes each pixel individually. The transformation follows a weighted luminance formula that prioritizes different color components based on their contribution to perceived brightness:

$$Gray = 0.3R + 0.59G + 0.11B$$

This formula reflects human visual perception, where green has the highest influence on brightness, followed by red and blue. The algorithm iterates through each pixel of the input image, applies this formula, and generates a new single-channel grayscale image optimized for text extraction.

To enhance efficiency and maintain image quality, the grayscale conversion function is optimized using the `SixLabors.ImageSharp` library. The `Apply` method ensures precise pixel-by-pixel processing using the following approach:

```
// Get the original pixel color at (x, y).
Rgba32 pixelColor = original[x, y];

// Compute grayscale intensity using the
// luminance formula.
byte gray = (byte)(0.3 * pixelColor.R + 0.59 *
    pixelColor.G + 0.11 * pixelColor.B);

// Assign the grayscale value to the new image.
grayscaleImage[x, y] = new L8(gray);
```

Grayscale conversion significantly reduces unnecessary color variations, making text stand out against the background. By providing a cleaner input for subsequent preprocessing techniques, this step ensures improved OCR accuracy and enhances the overall robustness of text recognition applications.

**2) Global Thresholding:** Global thresholding is a foundational technique in image preprocessing, especially useful for converting grayscale or color images into binary format for Optical Character Recognition (OCR) systems. The process involves applying a single, fixed threshold value across the entire image. Pixels with grayscale intensity values greater than or equal to this threshold are set to white (255), while those below are set to black (0). This transformation produces a high-contrast binary image, simplifying the task of identifying text against a background.

This method is particularly effective for images with uniform lighting conditions, where pixel intensity distribution is relatively consistent. However, global thresholding may underperform in scenarios involving uneven illumination, shadows, or variable backgrounds, as it lacks adaptive capabilities.

In our project, we implemented global thresholding using the `ImageSharp` library. The core algorithm is encapsulated in a dedicated `GlobalThresholding` class, which receives a threshold value at instantiation.

```
byte value = grayscale >= _threshold ? (byte)255
    : (byte)0;
thresholdedImage[x, y] = new Rgba32(value, value,
    value, 255);
```

A commonly used threshold value falls within the range of 80 to 150, depending on the contrast and clarity of the input image. Multiple instances of the `GlobalThresholding` class were initialized with varying thresholds ranging from 80 to 150, which are commonly effective for ID card and document images.

**3) Adaptive Thresholding:** Unlike global thresholding, which uses a fixed cutoff value for all pixels, adaptive thresholding dynamically computes the binarization threshold for each pixel based on the local intensity distribution in its surrounding neighborhood. This makes it particularly effective for images exhibiting non-uniform illumination, shadows, or uneven exposure, where a single global threshold would be inadequate.

The threshold for each region is computed using the local mean intensity:

```
byte localThreshold = (byte)((sum / blockArea) -
    c);
```

Adaptive thresholding enhances OCR accuracy by preserving fine details in the text while minimizing noise interference. This technique is particularly useful for scanned documents, handwritten text, and images captured in challenging lighting conditions. By selecting the appropriate thresholding method based on the input image characteristics, OCR performance can be significantly improved. We configured the adaptive thresholding with a `blockSize` of 11 and a `c` value of 5.0, which provided consistently clean binarizations across a diverse set of input images. A `blockSize` of 11 ensures that each pixel's threshold is calculated using a sufficiently large local neighborhood, capturing enough contextual detail without being overly influenced by small noise. If the block size were smaller, it might react too sensitively to noise, while a larger block size could blur important text features. The `C` value of 5.0 helps adjust the threshold downward slightly, preventing areas with subtle fading or shadows from being incorrectly classified as background. This value helps fine-tune the threshold to be more forgiving of minor imperfections like faded ink, which is especially important for documents with variable lighting or uneven print quality. The combination of a moderate block size and an appropriate `C` value ensures robust binarization across a diverse set of images, especially those with varying contrast or printing flaws.

**4) Saturation Adjustment:** Saturation adjustment is an essential preprocessing technique that plays a significant role in enhancing the clarity of text in images, especially in cases where the text is faded, low-contrast, or noisy. By modifying the intensity of the colors in an image, saturation adjustment increases the contrast between the text and its background. This improvement in contrast is particularly beneficial for Optical Character Recognition (OCR) accuracy, as it helps make the text more distinguishable from the background. Such enhancement is crucial before other preprocessing steps

like binarization and thresholding, as it ensures better text segmentation, leading to improved OCR performance.

In the context of this project, saturation adjustment is implemented using the SaturationAdjustment class. The core functionality of this class lies in its Apply method, which processes the image by adjusting the color intensity of each pixel based on a given saturation factor. The transformation is performed efficiently with the SixLabors.ImageSharp library, which provides tools for seamless manipulation of the image's pixels while maintaining the image structure.

The first step in the implementation involves cloning the original image. This ensures that the integrity of the original image is preserved, and any modifications are applied to a separate copy, preventing changes to the source image. A color transformation matrix is then applied, which adjusts the RGB (Red, Green, Blue) channels of each pixel based on the saturation factor. This modification alters the intensity of each color in the image, effectively changing its saturation. After applying the transformation, the new pixel values are clamped within the valid RGB range of 0 to 255. This clamping process ensures that the pixel values remain valid and prevents color distortion that may occur if the values exceed the allowed range. Finally, the alpha (transparency) channel is left unchanged, preserving the image's original transparency and overall structure.

Here is a code snippet that illustrates the process of pixel manipulation involved in saturation adjustment:

```
Rgba32 pixel = adjustedImage[x, y];

// Retrieve the original pixel color
Rgba32 pixel = adjustedImage[x, y];

// Apply the color transformation based on the
// saturation factor
float newRed = pixel.R * colorMatrixElements
    [0][0] + pixel.G * colorMatrixElements[1][0]
    + pixel.B * colorMatrixElements[2][0];
float newGreen = pixel.R * colorMatrixElements
    [0][1] + pixel.G * colorMatrixElements[1][1]
    + pixel.B * colorMatrixElements[2][1];
float newBlue = pixel.R * colorMatrixElements
    [0][2] + pixel.G * colorMatrixElements[1][2]
    + pixel.B * colorMatrixElements[2][2];

// Clamp pixel values to maintain valid RGB range
adjustedImage[x, y] = new Rgba32(
    (byte)Math.Clamp(newRed, 0, 255),
    (byte)Math.Clamp(newGreen, 0, 255),
    (byte)Math.Clamp(newBlue, 0, 255),
    pixel.A // Preserve alpha transparency
);
```

This saturation adjustment process effectively enhances the image's color contrast, making the text more visible and improving OCR accuracy. In scanned documents where the text may have faded or blended into the background, adjusting the saturation helps separate the text from the background, enabling more accurate text recognition. This preprocessing step is essential for ensuring reliable OCR results, making it a critical component in the image preprocessing pipeline for text extraction tasks.

*5) Deskewing (Skew Correction):* Skewed text in scanned or captured images can negatively impact OCR accuracy by distorting character alignment, making text recognition challenging. Deskewing, also known as skew correction, is a preprocessing step that detects and corrects text orientation to ensure proper alignment before OCR processing.

We implemented a deskewing module as part of the preprocessing pipeline using the ImageSharp library in C#. The process begins with the conversion of the input image to grayscale, which reduces computational complexity and focuses the analysis on luminance variations. This is followed by edge detection, performed using a Sobel filter approximation, to emphasize structural lines within the document—particularly those corresponding to text baselines and lines of printed content. After edge detection, the system calculates the skew angle, which represents the degree of deviation from the horizontal alignment.

```
double skewAngle = CalculateSkewAngle(DetectEdges
    (inputImage));
```

After determining the skew angle, the image is rotated in the opposite direction to align the text horizontally, ensuring that the OCR system can accurately interpret characters without distortion.

```
return RotateImage(inputImage, -skewAngle);
```

By aligning text properly, deskewing significantly enhances OCR performance, improving recognition accuracy, readability, and overall efficiency in text extraction. Integrating this step into the preprocessing pipeline ensures more reliable OCR results.

## B. OCR Processing Using the Terrasect SDK

After preprocessing the images, text extraction is performed using the Tesseract SDK. This open-source OCR engine is integrated into a C# console application to automate text recognition and support batch processing. The implementation ensures efficient and scalable text extraction from processed images.

Within the application, a dedicated function handles text extraction from images, leveraging the capabilities of Tesseract. The function processes the image, applies the specified OCR model, and stores the extracted text for further analysis.

```
string extractedText = TesseractProcessor.
    ExtractTextFromImage(processedImage,
    modelName, createdFilePath, fileWriter);
```

Once the text is extracted, it is stored for evaluation. The results from different preprocessing techniques are compared to determine the most effective method in enhancing OCR accuracy. This approach ensures that text recognition is optimized, leading to improved reliability in real-world applications.

### C. Performance Evaluation Metrics

To evaluate the effectiveness of each preprocessing method, three key performance indicators were analyzed:

The first metric, cosine similarity, measures how closely the extracted text from one preprocessing technique matches other preprocessing technique. A higher similarity score indicates better text extraction accuracy, making this metric crucial for assessing OCR performance.

The second metric, Processing Time, tracks the duration taken to process an image through each preprocessing method. A high-precision timer records the execution time from start to completion, ensuring a fair comparison between different approaches. Faster execution times indicate more efficient preprocessing, which is especially important for large-scale OCR applications where processing speed impacts overall system performance.

The third metric, Memory Consumption, monitors the amount of system memory used during image preprocessing. This metric is essential for evaluating computational efficiency, particularly for handling high-resolution images or large datasets. Lower memory usage is preferred, as it enables processing on resource-constrained environments without excessive system overhead.

### D. Comparative Analysis and Model Selection

A weighted scoring approach is used to determine the optimal preprocessing model by balancing accuracy, processing speed, and memory efficiency. The cosine similarity metric contributes the highest weight, as it directly affects OCR output quality. Processing speed is also prioritized, as faster methods are preferable in real-time or large-scale applications. Lastly, memory usage is considered to ensure that the chosen method remains computationally efficient.

By assigning appropriate weightages to each criterion, the overall performance of different preprocessing techniques can be ranked, allowing for an informed selection of the best approach. This ranking system ensures an optimal trade-off between accuracy, efficiency, and computational overhead, leading to improved OCR results while maintaining a balance between performance and resource utilization.

## IV. IMPLEMENTATION

The OCR pipeline integrates embedding generation and cosine similarity calculation to evaluate the performance of various preprocessing models. The implementation workflow is illustrated in Fig. 1:

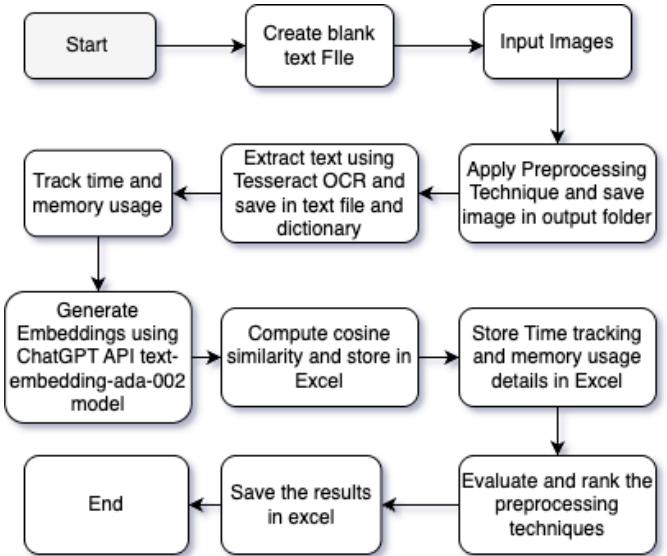


Fig. 1: Workflow Diagram

### A. Initialization and Service Setup

The pipeline begins by initializing key services and loading configuration settings. Logging is enabled to capture events throughout the process. The EmbeddingGeneratorService generates embeddings for extracted text using OpenAI's API. Performance tracking is handled by the ProcessingTimeTracker and ProcessingMemoryTracker, measuring the time and memory usage during preprocessing and text extraction. Cosine similarity calculations are performed to assess the similarity between text outputs from different preprocessing models. Configuration settings, including file paths for input images, output directories, and API keys, are loaded at the start.

### B. Image Preprocessing and OCR Execution

Before extracting text, images undergo several preprocessing steps to optimize them for OCR:

- **Resizing:** Ensures uniform dimensions for consistent OCR performance.
- **Grayscale Conversion:** Enhances contrast and reduces noise.
- **Adaptive Thresholding:** Dynamically adjusts threshold values for improved text visibility under varying lighting conditions.
- **Saturation Adjustments:** Applied at different levels (0.6, 0.9, 1.2, 1.6, 2.0, 2.5) to improve text clarity.
- **Deskewing:** Corrects misaligned text for better OCR accuracy.
- **Global Thresholding:** Converts images to grayscale to enhance contrast.

### C. Text Extraction with Tesseract OCR

After preprocessing, text extraction is performed using the Tesseract OCR engine. The processed image is loaded into memory, converted into a PNG stream using the SaveAsPng method, and passed to Tesseract. The image is loaded as a Pix

object, and Tesseract recognizes the text. The extracted text is retrieved using `page.GetText()` and saved for future use. To maintain organization across multiple preprocessing models, extracted text is labeled with metadata indicating the image source, processing type, or OCR model used. The labeled text is then stored using the `FileWriter` class.

#### D. Performance Tracking

The system records processing time and memory usage for each preprocessing model:

- **Processing Time:** The `ProcessingTimeTracker` measures the time taken for each transformation and OCR execution in milliseconds.
- **Memory Usage:** The `ProcessingMemoryTracker` records the memory consumed during processing using the `MeasureMemoryUsage` method.

#### E. Generating Embeddings and Cosine Similarity Calculation

After text extraction, embeddings are generated, and cosine similarity is calculated to compare results across different preprocessing models. The extracted text is converted into numerical vectors using embeddings generated via OpenAI's API. The `EmbeddingGeneratorService` ensures each preprocessing model—such as grayscale conversion or thresholding—produces a corresponding embedding stored as a float array.

Cosine similarity is used to measure textual similarity. The similarity score ranges from 1 (identical text) to 0 (completely dissimilar text). The `CosineSimilarityCalculator` computes this metric using:

$$\text{cosine similarity} = \frac{A \cdot B}{\|A\| \|B\|}$$

where  $A \cdot B$  is the dot product of two vectors, and  $\|A\|$  and  $\|B\|$  are their magnitudes.

The similarity scores are stored in an Excel report, facilitated by the NPOI library. OpenAI embedding generation requests are sent via `HttpClient`, ensuring an automated process.

#### F. Evaluating Preprocessing Models

The extracted text and similarity results are analyzed to identify the most effective preprocessing models. Cosine similarity helps determine text consistency, while processing time and memory usage help in selecting efficient techniques. The `PreprocessingModelEvaluator` ranks models based on a weighted score:

- **Cosine Similarity:** 50% weight
- **Processing Time:** 30% weight
- **Memory Usage:** 20% weight

The best-performing models are identified and stored for further use.

#### G. Output and Data Export

The extracted text, along with metadata such as file names and timestamps, is stored and exported using the NPOI library. Key operations, errors, and warnings are logged for debugging and auditing.

## V. RESULTS

**1) Performance Evaluation:** The effectiveness of text extraction is assessed through a structured evaluation of key performance metrics.

**Impact of Preprocessing on Text Extraction:** The images below illustrate the impact of different preprocessing techniques on input image text.



Fig. 2: Input image of an ID Card

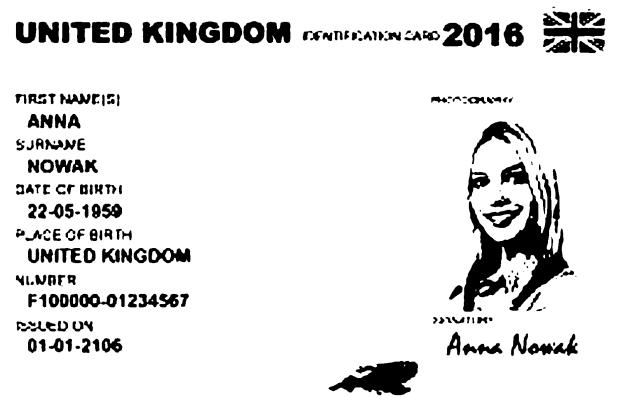


Fig. 3: After Global Thresholding at 150



Fig. 4: After Saturation adjustment at 0.6



Fig. 5: After Saturation adjustment at 2.5



Fig. 6: After Deskewing Image



Fig. 7: After Adaptive Thresholding

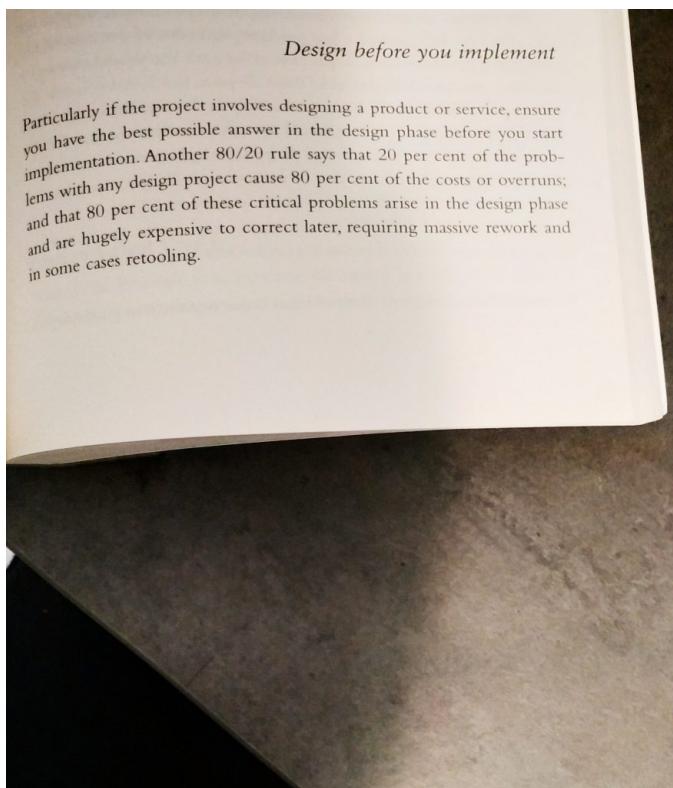


Fig. 8: Input Image of Book Page with Lighting variation

*Design before you implement*

Particularly if the project involves designing a product or service, ensure you have the best possible answer in the design phase before you start implementation. Another 80/20 rule says that 20 per cent of the problems with any design project cause 80 per cent of the costs or overruns; and that 80 per cent of these critical problems arise in the design phase and are hugely expensive to correct later, requiring massive rework and in some cases retooling.

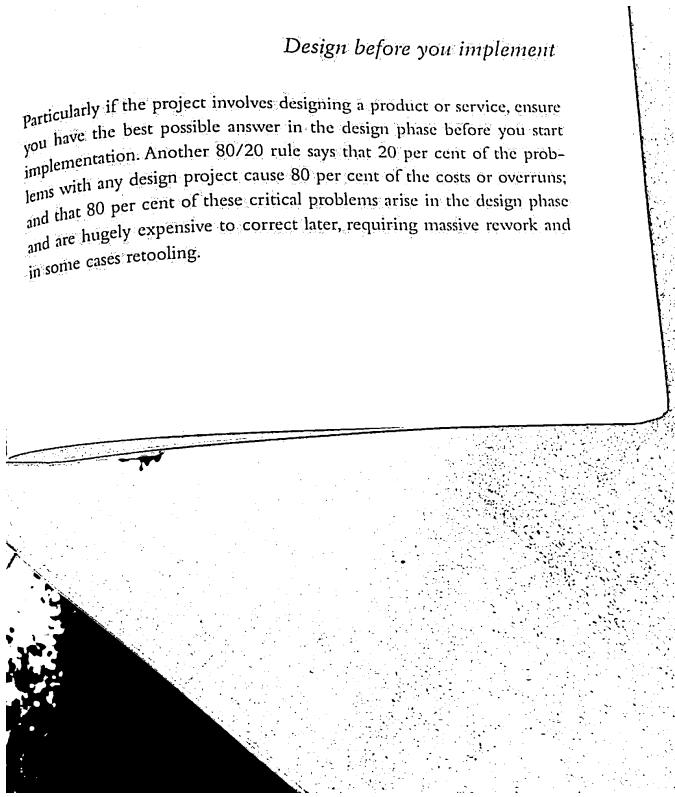


Fig. 9: After Adaptive Thresholding

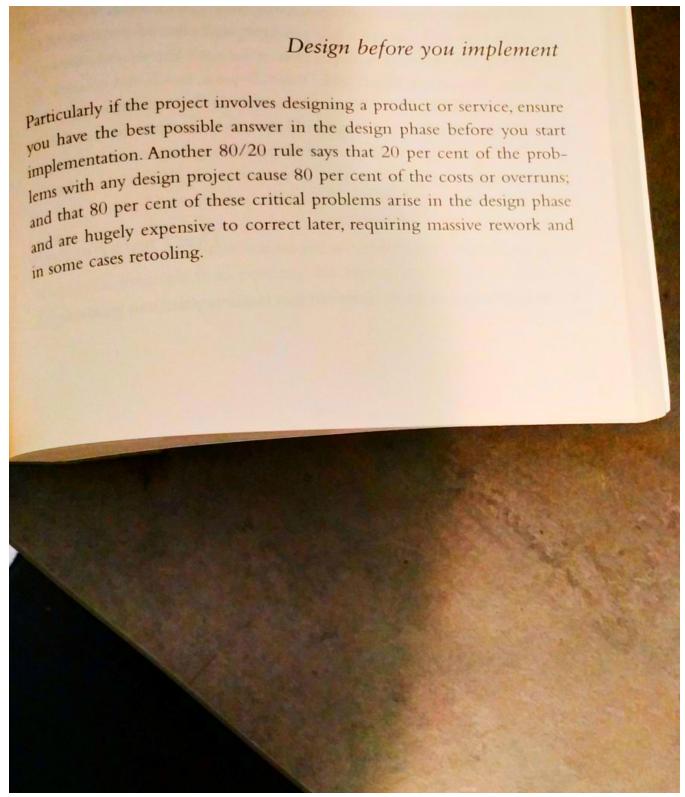


Fig. 11: After Saturation Adjusted at 2.5

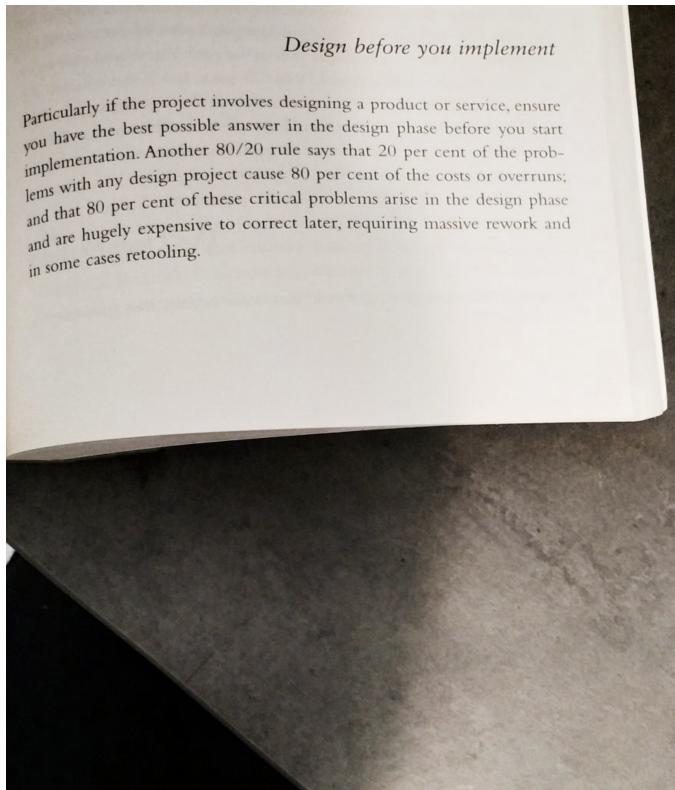


Fig. 10: After Saturation Adjustedment at 0.6

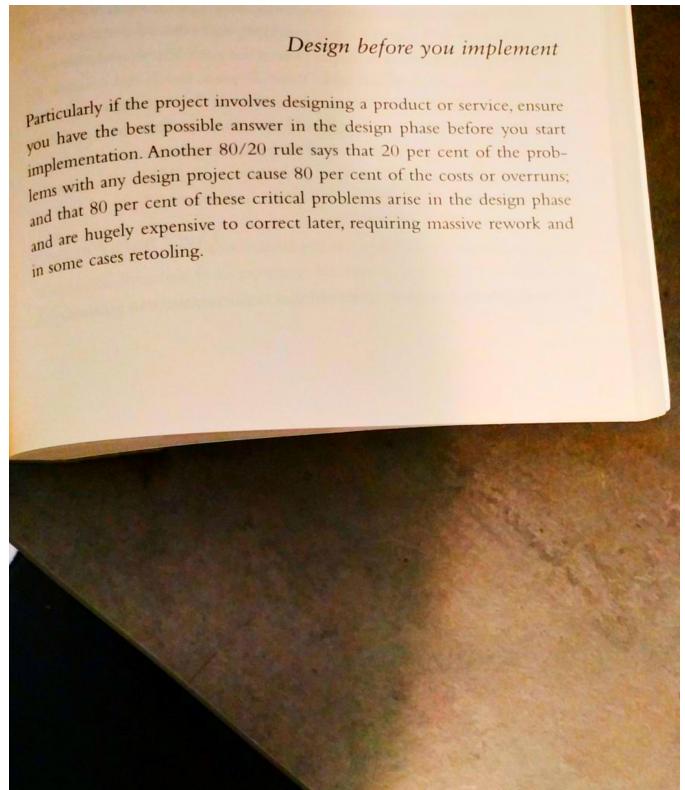


Fig. 12: After Saturation Adjusted at 2.5

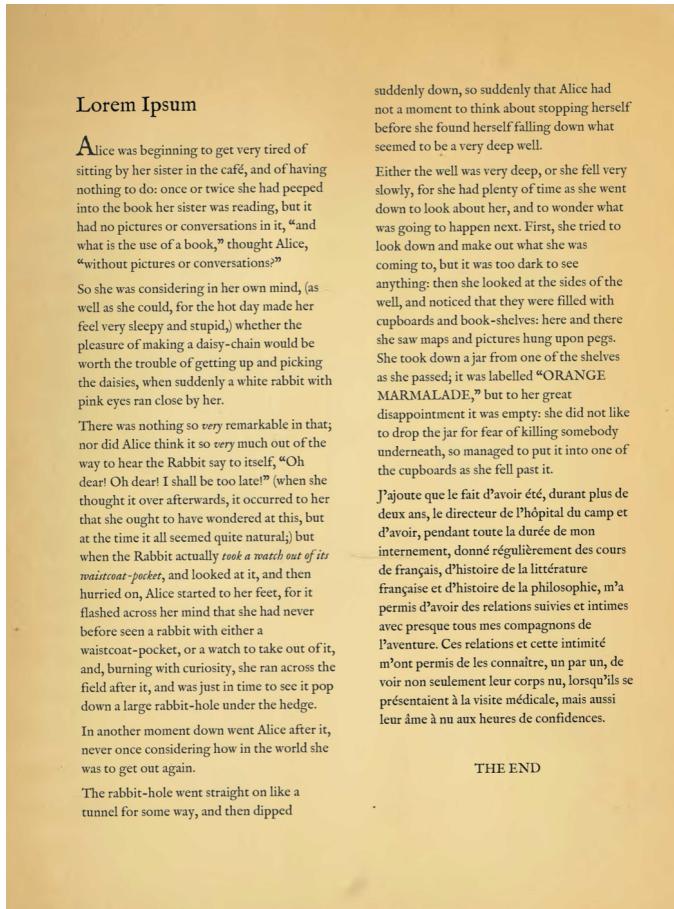


Fig. 13: Input image



Fig. 14: After converting to Grayscale

Grayscale conversion eliminates unnecessary color details, improving text visibility. Thresholding converts images into pure black and white, aiding OCR but potentially removing finer details.



Fig. 15: After global thresholding at 80

### **Lore Ipsum**

Alice was beginning to get very tired of sitting by her sister in the café, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, "and what is the use of a book?" thought Alice, "without pictures or conversations?"

So she was considering in her own mind, (as well as she could, for the hot day made her feel very sleepy and stupid,) whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a white rabbit with pink eyes ran close by her.

There was nothing so *very* remarkable in that; nor did Alice think it so *very* much out of the way to hear the Rabbit say to itself, "Oh dear! Oh dear! I shall be too late!" (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually *took a watch out of its waistcoat-pocket*, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket, or a watch to take out of it, and, burning with curiosity, she ran across the field after it, and was just in time to see it pop down a large rabbit-hole under the hedge.

In another moment down went Alice after it, never once considering how in the world she was to get out again.

The rabbit-hole went straight on like a tunnel for some way, and then dipped

suddenly down, so suddenly that Alice had not a moment to think about stopping herself before she found herself falling down what seemed to be a very deep well.

Either the well was very deep, or she fell very slowly, for she had plenty of time as she went down to look about her, and to wonder what was going to happen next. First, she tried to look down and make out what she was coming to, but it was too dark to see anything; then she looked at the sides of the well, and noticed that they were filled with cupboards and book-shelves: here and there she saw maps and pictures hung upon pegs. She took down a jar from one of the shelves as she passed; it was labelled "ORANGE MARMALADE," but to her great disappointment it was empty: she did not like to drop the jar for fear of killing somebody underneath, so managed to put it into one of the cupboards as she fell past.

J'ajoute que le fait d'avoir été, durant plus de deux ans, le directeur de l'hôpital du camp et d'avoir, pendant toute la durée de mon internement, donné régulièrement des cours de français, d'histoire de la littérature française et d'histoire de la philosophie, m'a permis d'avoir des relations suivies et intimes avec presque tous mes compagnons de l'aventure. Ces relations et cette intimité m'ont permis de les connaître, un par un, de voir non seulement leur corps nu, lorsqu'ils se présentaient à la visite médicale, mais aussi leur âme à nu aux heures de confidences.

THE END

Fig. 16: After global thresholding at 150

Saturation adjustment enhances contrast between the text and background, making the content more distinguishable.

### **Lore Ipsum**

Alice was beginning to get very tired of sitting by her sister in the café, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, "and what is the use of a book?" thought Alice, "without pictures or conversations?"

So she was considering in her own mind, (as well as she could, for the hot day made her feel very sleepy and stupid,) whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a white rabbit with pink eyes ran close by her.

There was nothing so *very* remarkable in that; nor did Alice think it so *very* much out of the way to hear the Rabbit say to itself, "Oh dear! Oh dear! I shall be too late!" (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually *took a watch out of its waistcoat-pocket*, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a waistcoat-pocket, or a watch to take out of it, and, burning with curiosity, she ran across the field after it, and was just in time to see it pop down a large rabbit-hole under the hedge.

In another moment down went Alice after it, never once considering how in the world she was to get out again.

The rabbit-hole went straight on like a tunnel for some way, and then dipped

suddenly down, so suddenly that Alice had not a moment to think about stopping herself before she found herself falling down what seemed to be a very deep well.

Either the well was very deep, or she fell very slowly, for she had plenty of time as she went down to look about her, and to wonder what was going to happen next. First, she tried to look down and make out what she was coming to, but it was too dark to see anything; then she looked at the sides of the well, and noticed that they were filled with cupboards and book-shelves: here and there she saw maps and pictures hung upon pegs. She took down a jar from one of the shelves as she passed; it was labelled "ORANGE MARMALADE," but to her great

disappointment it was empty: she did not like to drop the jar for fear of killing somebody underneath, so managed to put it into one of the cupboards as she fell past.

J'ajoute que le fait d'avoir été, durant plus de deux ans, le directeur de l'hôpital du camp et d'avoir, pendant toute la durée de mon internement, donné régulièrement des cours de français, d'histoire de la littérature

française et d'histoire de la philosophie, m'a

permis d'avoir des relations suivies et intimes

avec presque tous mes compagnons de

l'aventure. Ces relations et cette intimité

m'ont permis de les connaître, un par un,

de voir non seulement leur corps nu, lorsqu'ils se

présentaient à la visite médicale, mais aussi

leur âme à nu aux heures de confidences.

THE END

Fig. 17: Saturation Adjustment at 0.6

### **Lore Ipsum**

Alice was beginning to get very tired of sitting by her sister in the café, and of having nothing to do: once or twice she had peeped into the book her sister was reading, but it had no pictures or conversations in it, "and what is the use of a book?" thought Alice, "without pictures or conversations?"

So she was considering in her own mind, (as well as she could, for the hot day made her feel very sleepy and stupid,) whether the pleasure of making a daisy-chain would be worth the trouble of getting up and picking the daisies, when suddenly a white rabbit with pink eyes ran close by her.

There was nothing so *very* remarkable in that; nor did Alice think it so *very* much out of the way to hear the Rabbit say to itself, "Oh dear! Oh dear! I shall be too late!" (when she thought it over afterwards, it occurred to her that she ought to have wondered at this, but at the time it all seemed quite natural); but when the Rabbit actually *took a watch out of its waistcoat-pocket*, and looked at it, and then hurried on, Alice started to her feet, for it flashed across her mind that she had never before seen a rabbit with either a

waistcoat-pocket, or a watch to take out of it,

and, burning with curiosity, she ran across the field after it, and was just in time to see it pop down a large rabbit-hole under the hedge.

In another moment down went Alice after it, never once considering how in the world she was to get out again.

The rabbit-hole went straight on like a tunnel for some way, and then dipped

suddenly down, so suddenly that Alice had not a moment to think about stopping herself before she found herself falling down what seemed to be a very deep well.

Either the well was very deep, or she fell very slowly, for she had plenty of time as she went down to look about her, and to wonder what was going to happen next. First, she tried to look down and make out what she was coming to, but it was too dark to see anything; then she looked at the sides of the well, and noticed that they were filled with cupboards and book-shelves: here and there she saw maps and pictures hung upon pegs. She took down a jar from one of the shelves as she passed; it was labelled "ORANGE MARMALADE," but to her great disappointment it was empty: she did not like to drop the jar for fear of killing somebody underneath, so managed to put it into one of the cupboards as she fell past.

J'ajoute que le fait d'avoir été, durant plus de deux ans, le directeur de l'hôpital du camp et d'avoir, pendant toute la durée de mon

internement, donné régulièrement des cours de français, d'histoire de la littérature

française et d'histoire de la philosophie, m'a

permis d'avoir des relations suivies et intimes

avec presque tous mes compagnons de

l'aventure. Ces relations et cette intimité

m'ont permis de les connaître, un par un,

de voir non seulement leur corps nu, lorsqu'ils se

présentaient à la visite médicale, mais aussi

leur âme à nu aux heures de confidences.

THE END

Fig. 18: Saturation adjustment at 2.5

**Text Accuracy (Cosine Similarity):** The accuracy of the extracted text is evaluated by comparing it with the original content using cosine similarity. This metric measures how closely the OCR-processed text aligns with the reference text. A higher similarity score signifies better text recognition and minimal deviation from the original text.

Model	Grayscale	GlobalThr	GlobalThr	GlobalThr	GlobalThr	GlobalThr	Saturation	Deskewed	AdaptiveThreshold								
Grayscale	1	0.9262	0.9357	0.9509	0.9321	0.9214	0.9049	0.9180	0.9079	0.9092	0.9318	0.9239	0.9204	0.9148	0.8817		
GlobalThreshold_80	0.92562	1	0.9594	0.9242	0.9188	0.9318	0.8899	0.8723	0.8828	0.8992	0.8971	0.8804	0.9225	0.85245			
GlobalThreshold_100	0.93571	0.95949	1	0.94309	0.9118	0.92468	0.91076	0.91076	0.91042	0.91037	0.91848	0.91775	0.91075	0.9318	0.87532		
GlobalThreshold_110	0.95059	0.9242	0.94308	1	0.92099	0.90021	0.88367	0.88367	0.88313	0.88076	0.89362	0.89956	0.90201	0.90922	0.86515		
GlobalThreshold_120	0.93271	0.91888	0.93118	0.92099	1	0.95764	0.94456	0.95015	0.94604	0.93695	0.94881	0.94335	0.94839	0.94461	0.93175		
GlobalThreshold_140	0.92114	0.88318	0.92468	0.9102	0.95764	1	0.98076	0.97651	0.98222	0.97407	0.97455	0.97262	0.97545	0.95961	0.94466		
GlobalThreshold_150	0.9049	0.88504	0.91076	0.88569	0.94586	0.98045	1	0.97693	0.98016	0.96528	0.96848	0.97086	0.97243	0.94279	0.93606		
SaturationAdjusted_1.2	0.91801	0.8899	0.914	0.89367	0.90505	0.97651	0.97603	1	0.98277	0.97168	0.97166	0.97122	0.96633	0.94264			
SaturationAdjusted_0.6	0.90791	0.87723	0.90421	0.8813	0.94604	0.98222	0.9806	0.98277	1	0.97771	0.97625	0.97284	0.97963	0.93142	0.93401		
SaturationAdjusted_0.9	0.9092	0.88288	0.91037	0.88076	0.93695	0.97407	0.96528	0.97168	0.97771	1	0.96812	0.96111	0.95667	0.94967	0.91763		
SaturationAdjusted_1.6	0.91913	0.88992	0.91848	0.90562	0.94881	0.97445	0.96848	0.97166	0.97625	0.96812	1	0.98686	0.97777	0.95556	0.92607		
SaturationAdjusted_2	0.92329	0.89771	0.91775	0.89956	0.94335	0.97825	0.97088	0.97122	0.97824	0.96111	0.98866	1	0.96975	0.94765	0.92117		
SaturationAdjusted_2.5	0.92014	0.88404	0.91705	0.90021	0.94839	0.97545	0.97243	0.96633	0.97386	0.95867	0.97777	0.96975	1	0.94394	0.92827		
Deskewed	0.94148	0.92225	0.93811	0.90922	0.94461	0.95968	0.94297	0.952	0.95142	0.94967	0.95556	0.94705	0.94394	1	0.91903		
AdaptiveThreshold	0.88172	0.85245	0.87532	0.86515	0.91375	0.94463	0.93605	0.92464	0.93404	0.91763	0.92607	0.92177	0.92827	0.91693	1		

Fig. 19: Cosine matrix result

**Processing Time:** The speed at which an image is processed plays a crucial role in determining the efficiency of a preprocessing technique. The time taken to execute each method is recorded, facilitating a comparative analysis. Faster processing time indicates a more optimized approach, making it suitable for large-scale applications where efficiency is a priority.

**Memory Consumption:** The computational efficiency of a preprocessing method is evaluated by monitoring its memory usage. Reducing memory consumption is essential for han-

dling large datasets efficiently without overloading system resources. A lower memory footprint ensures smooth processing while maintaining accuracy.

Model	Average Time (s)	Average Memory (MB)
Grayscale	2024.865767	0.013807954
GlobalThreshold_80	1676.491633	0.002740924
GlobalThreshold_100	1701.283533	0.001915953
GlobalThreshold_110	1733.68175	0.00275027
GlobalThreshold_120	1655.744083	0.00275027
GlobalThreshold_140	1670.999433	0.00275027
GlobalThreshold_150	1693.374417	0.002975761
SaturationAdjusted_1.2	2127.863867	0.001956601
SaturationAdjusted_0.6	2134.148967	0.00275027
SaturationAdjusted_0.9	2123.34895	0.002707142
SaturationAdjusted_1.6	2145.41245	0.002723567
SaturationAdjusted_2	2085.000333	0.002681414
SaturationAdjusted_2.5	2074.20855	0.00275027
Deskewed	1199.244817	2.110389561
AdaptiveThreshold	4515.6431	0.002624088

Fig. 20: Time tracking and memory usage ranking

By analyzing these performance metrics, an optimal preprocessing strategy can be identified, balancing accuracy, speed, and resource efficiency. A scoring mechanism is employed to rank the different preprocessing models, integrating cosine similarity, execution time, and memory efficiency to determine the most effective approach for enhancing OCR performance.

Rank	Model	Final Score
1	Grayscale	0.86417
2	GlobalThreshold_110	0.86002
3	GlobalThreshold_120	0.85609
4	GlobalThreshold_100	0.85465
5	GlobalThreshold_80	0.85117
6	GlobalThreshold_140	0.84929
7	GlobalThreshold_150	0.83967
8	SaturationAdjusted_1.6	0.82312
9	SaturationAdjusted_2	0.82287
10	SaturationAdjusted_2.5	0.82201
11	SaturationAdjusted_1.2	0.81745
12	SaturationAdjusted_0.9	0.81328
13	SaturationAdjusted_0.6	0.81191
14	Deskewed	0.69107
15	AdaptiveThreshold	0.64061

Fig. 21: Final ranking

## VI. CONCLUSION

This paper presents an enhanced Optical Character Recognition (OCR) system that integrates advanced image preprocessing techniques with Tesseract OCR, cosine similarity analysis, and embedding-based text representation using OpenAI's text-embedding-ada-002 model. The goal of this research was to identify the most effective pre-processing techniques for improving OCR accuracy while maintaining computational efficiency in terms of time and memory usage. Through extensive experimentation, we evaluated multiple preprocessing techniques, including grayscale conversion, adaptive thresholding, global thresholding, saturation adjustment,

and deskewing. Each of these techniques was assessed based on its impact on OCR accuracy, processing time, and memory consumption. Our findings indicate that a combination of grayscale conversion and Global thresholding at threshold 110 yields the best balance between accuracy and computational efficiency. Grayscale conversion simplifies the image by reducing color variations, thereby improving character differentiation. Global thresholding further enhances the contrast between text and background, making the extracted text more distinguishable to the OCR engine.

In addition to traditional OCR processing, this study incorporated text similarity analysis using the Cosine Similarity Matrix to improve recognition reliability. Furthermore, we explored the benefits of embedding generation using OpenAI's text-embedding-ada-002 model. This model generates numerical representations of text, capturing semantic meaning and contextual relationships. By leveraging text embeddings, the OCR system enhances its capability to perform text classification, clustering, and search tasks. Our results demonstrate that embedding-based text analysis significantly improves document organization and retrieval accuracy, making it a valuable addition to modern OCR pipelines.

One of the unique contributions of this research is the evaluation of computational efficiency, which is often overlooked in OCR studies. We recorded the time and memory usage for each pre-processing technique and observed that while some methods improve accuracy, they also introduce additional computational overhead. The balance between accuracy and efficiency is crucial for real-world applications where large-scale document processing is required. Our findings suggest that grayscale conversion and adaptive thresholding provide the optimal trade-off, ensuring high accuracy without excessive resource consumption. Overall, this study reinforces the importance of pre-processing techniques in enhancing OCR performance. By integrating Tesseract OCR with cosine similarity and embedding-based text representation, we have demonstrated a novel approach that improves text recognition accuracy while optimizing computational resources. The insights from this research can be applied to various domains, including automated data entry, historical document preservation, and intelligent document retrieval systems.

For future work, we propose further exploration of deep learning-based image enhancement techniques, such as Generative Adversarial Networks (GANs) for document restoration. Additionally, integrating multimodal learning approaches that combine textual and visual information could further refine OCR accuracy. The continuous advancement of AI-driven text recognition methods presents exciting opportunities for future research and development in this field.

## VII. LIMITATIONS AND FUTURE WORK

While this research has compared various preprocessing techniques, there remain several opportunities for future exploration. One promising direction is the integration of deep learning-based image enhancement techniques, such as Generative Adversarial Networks (GANs) or Convolutional Neural

Networks (CNNs), to further improve image quality before text extraction. These advanced models can effectively remove noise, enhance text contrast, and reconstruct degraded portions of images, leading to improved OCR results. Additionally, future work can focus on the development of adaptive pre-processing pipelines that dynamically select the best enhancement techniques based on image characteristics. Instead of applying a fixed set of pre-processing steps, machine learning models can be trained to determine the optimal transformation techniques for each input image, maximizing accuracy while minimizing unnecessary computational costs. Another potential research avenue is the exploration of multi-modal learning, which combines OCR with additional contextual information, such as document structure, handwriting recognition, and natural language processing (NLP). Integrating OCR with NLP-based correction models can further improve text recognition by refining misclassified words based on semantic understanding and contextual analysis. Moreover, real-world OCR applications often involve multi-language documents, requiring robust language identification and character recognition capabilities. Future research can investigate the effectiveness of transformer-based models, such as multilingual BERT, in improving OCR for diverse languages and scripts. Finally, extending this study to large-scale datasets and real-time applications would provide valuable insights into the practical deployment of OCR systems in industries such as finance, healthcare, and legal documentation. Evaluating the scalability and adaptability of the proposed approach in high-volume text processing scenarios will be crucial for enhancing OCR adoption in enterprise applications.

### VIII. ACKNOWLEDGEMENT

We sincerely appreciate the invaluable guidance and insightful feedback provided by Professor Dr. Damir Dobric throughout this project. His support has been instrumental in shaping our research. Additionally, we extend our gratitude to the tutors for their assistance in clarifying doubts and offering their valuable insights.

### REFERENCES

- [1] R. Smith, "An overview of the Tesseract OCR engine," in *Proc. 9th Int. Conf. Document Anal. Recognit. (ICDAR)*, Curitiba, Brazil, 2007, pp. 629–633. [Online]. Available: <https://ieeexplore.ieee.org/document/4376991>
- [2] W. Bieniecki and S. Grabowski, "Image preprocessing for improving OCR accuracy," in *Proc. Int. Conf. Perspective Technol. Methods MEMS Design (MEMSTECH)*, 2007. [Online]. Available: <https://ieeexplore.ieee.org/document/4283429>
- [3] W. Moon, S. Nengroo, T. Kim, J. Lee, S. Son, and D. Har, "Enhanced optical character recognition by optical sensor combined with BERT and cosine similarity scoring (Student abstract)," in *Proc. AAAI Conf. Artif. Intell.*, 2025. [Online]. Available: <https://doi.org/10.1609/aaai.v38i21.30483>
- [4] OpenAI, "Introducing text-embedding-ada-002," OpenAI, 2025. [Online]. Available: <https://openai.com/blog/new-embedding-model>
- [5] N. Muda, N. K. Ismail, S. A. Abu Bakar, and J. M. Zain, "Optical character recognition by using template matching (alphabet)," 2009. [Online]. Available: [https://www.academia.edu/714194/Optical\\_Character\\_Recognition\\_By\\_Use\\_Template\\_Matching\\_Alphabet](https://www.academia.edu/714194/Optical_Character_Recognition_By_Use_Template_Matching_Alphabet)
- [6] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, pp. 436–444, May 2015. [Online]. Available: <https://www.nature.com/articles/nature14539>

- [7] K. Harraj and N. Raissouni, "A new pre-processing approach for improving OCR accuracy on documents captured by digital cameras under non-uniform illumination," *Zenodo*, Aug. 2018. [Online]. Available: <https://zenodo.org/record/1322323>
- [8] I. Singh, M. Colom, and K. Bontcheva, "A comparative analysis of OCR models on diverse datasets: Insights from Memes and Hiertext dataset," in *Proc. WACV VisionDocs Workshop*, 2025. [Online]. Available: [https://openaccess.thecvf.com/content/WACV2025W/VISIONDOCS/papers/Singh\\_A\\_Comparative\\_Analysis\\_of\\_OCR\\_Models\\_on\\_Diverse\\_Datasets\\_Insights\\_WACVW\\_2025\\_paper.pdf](https://openaccess.thecvf.com/content/WACV2025W/VISIONDOCS/papers/Singh_A_Comparative_Analysis_of_OCR_Models_on_Diverse_Datasets_Insights_WACVW_2025_paper.pdf)
- [9] N. Arica and F. T. Yarman-Vural, "An overview of character recognition focused on off-line handwriting," *IEEE Trans. Syst., Man, Cybern. C, Appl. Rev.*, vol. 31, no. 2, pp. 216–233, May 2001. [Online]. Available: [https://www.researchgate.net/publication/3421397\\_An\\_overview\\_of\\_character\\_recognition\\_focused\\_on\\_off-line\\_handwriting](https://www.researchgate.net/publication/3421397_An_overview_of_character_recognition_focused_on_off-line_handwriting)
- [10] M. Li *et al.*, "TrOCR: Transformer-based optical character recognition with pre-trained models," *arXiv preprint*, arXiv:2109.10282, 2021. [Online]. Available: <https://arxiv.org/abs/2109.10282>