

Enhanced Scene Text Extraction through “Texture Analysis and Deep Convolutional Networks”

Shilpi Rani

Lecturer Computer, Mahamaya Polytechnic of Information Technology, Amroha

ARTICLE INFO

Received: 29 Dec 2024

Revised: 15 Feb 2025

Accepted: 24 Feb 2025

ABSTRACT

The widespread use of portable cameras and advancements in visual computing have made extracting text from images captured in natural settings an increasingly important area of research. This capability can support various applications, including enhanced augmented reality experiences. Text extraction algorithms for complex scenes typically involve three main stages: (i) identifying and locating text regions, (ii) improving and isolating the text, and (iii) recognizing the characters using optical character recognition (OCR). However, this process is complicated by various challenges, such as inconsistent text sizes, different fonts and colors, diverse alignments, changes in lighting, and reflective surfaces. This paper reviews and categorizes current approaches, focusing primarily on the first two stages—text detection and segmentation—since the field of OCR is well-established and supported by reliable tools. Additionally, a publicly available image dataset is introduced to assist in evaluating and comparing methods for scene text extraction

Introduction: In recent decades, extracting text from visual media—such as images and videos—has become a topic of growing interest in the field of digital information management [1,2]. This technology has a wide range of applications, including automated video indexing, content summarization, searching, and retrieval [1,4]. A notable example is the Informedia project at Carnegie Mellon University, which utilizes textual content from newscasts and documentaries to enable comprehensive video search capabilities [5].

Objectives: The goal of this research is to enhance the precision and reliability of text extraction from images captured in natural scenes. This is achieved by combining texture analysis with deep convolutional neural networks (CNNs).

Methods: A wide range of research has addressed the issue of extracting text from natural scenes. These methods are Region-Based Methods ,Texture-Based Methods, Connected Component-Based Techniques, Edge Detection Approaches etc.

Results: The proposed method combining texture analysis with deep convolutional neural networks (CNNs) was evaluated using standard benchmark datasets . The results demonstrated significant improvements over traditional and deep learning-only approaches in various aspects:Improved Accuracy , Better Robustness in Complex Scenes, Higher Recall Rates, Reduced False Positives, Efficient Computation.

Conclusions: In this study, I have presented a robust method for text extraction from natural scene images by leveraging texture-based features and deep learning techniques. The integration of Histogram of Oriented Gradients (HOG) for texture feature extraction, followed by region validation using Convolutional Neural Networks (CNN), demonstrated significant effectiveness in identifying and isolating text from complex backgrounds. The proposed approach successfully addressed common challenges such as background clutter, varying illumination, and font diversity, which often hinder traditional OCR systems. Experimental results confirmed that my method enhances both detection accuracy and recognition reliability, especially in uncontrolled environments. Future work may focus on optimizing CNN architectures for real-time applications and expanding the model to support multilingual text recognition.

Keywords: Bounding Box, CNN, HOG, Identification, Localization, Local Binary Pattern, OCR, Segmentation Text Detection, Text Extraction, Texture-based,

INTRODUCTION

In recent decades, extracting text from visual media—such as images and videos—has become a topic of growing interest in the field of digital information management [1,2]. This technology has a wide range of applications, including automated video indexing, content summarization, searching, and retrieval [1,4]. A notable example is the Informedia project at Carnegie Mellon University, which utilizes textual content from newscasts and documentaries to enable comprehensive video search capabilities [5]. Text found within images often holds meaningful semantic information that aids in interpreting the visual content [2]. This text can typically be divided into two main types: **scene text** and **artificial text** [6]. Scene text refers to naturally occurring text present at the time the image was captured, such as on street signs, advertisements, or posters. In contrast, artificial text is superimposed on the image later, such as captions in television broadcasts or names shown during interviews.

Artificial text usually serves to clarify or describe the content of the image or video. On the other hand, scene text poses greater challenges due to its variation in appearance, including changes in orientation, size, style, and color[4]. As a result, extracting scene text accurately is more complex.

The process of automatically detecting and extracting text from natural scenes remains a highly challenging task. Key difficulties include:

1. **Diverse text characteristics:** text may differ in style, size, color, orientation, alignment, and position.
2. **Visual obstructions and quality degradation:** text can be blurred by motion or obscured by objects [7].
3. **Perspective distortions:** text captured in real-world scenes may appear skewed due to its placement in three-dimensional environments.

Text extraction from scene images presents several challenges due to the nature of the environments in which the text appears. Text can often be tilted or shaped based on the object's surface, and it can exhibit multiple orientations — horizontal, vertical, angled, or even mixed — as commonly seen on clothing or crumpled signage.

In recent years, extracting text from natural scenes has become a vital task in computer vision, allowing machines to interpret textual elements within everyday images. This capability plays a crucial role in various domains, such as autonomous vehicles, augmented reality, and accessibility tools. Nonetheless, the diverse and unpredictable nature of real-world environments—including varying lighting, complex backgrounds, and a broad array of font styles—presents significant obstacles for traditional text detection approaches.

To address these issues, combining texture analysis with deep convolutional neural networks (CNNs) has proven to be an effective strategy. Texture analysis helps in identifying intricate visual patterns that distinguish text from visually complex surroundings, while CNNs are adept at extracting layered features directly from image data. The integration of these methods results in more precise and resilient text extraction, particularly in scenes with noise or low visual contrast.

This study introduces a hybrid approach that fuses texture-based descriptors with deep learning techniques to advance the detection and recognition of scene text. By harnessing the strengths of both methodologies, the proposed system aspires to surpass existing solutions, offering improved flexibility and performance in varied real-world conditions.

OBJECTIVES

The objective of this research is to enhance the precision and reliability of text extraction from images captured in natural scenes. This is achieved by combining texture analysis with deep convolutional neural networks (CNNs). The proposed method is designed to accurately identify and extract text from challenging environments that include complex backgrounds, inconsistent lighting, and a wide range of font styles. Key components of the approach include:

Texture Analysis: Utilized to effectively differentiate text areas from the background by analyzing their distinctive texture properties, which aids in identifying text even in visually noisy or low-contrast settings.

Deep Convolutional Neural Networks (CNNs): Applied to learn features directly from data for detecting and recognizing text. This allows for robust performance across various scene conditions through end-to-end training.

Hybrid Feature Integration: Merging traditional handcrafted texture features with deep learning-derived features to strengthen the model's adaptability and boost the overall effectiveness of scene text extraction

METHODS

A wide range of research has addressed the issue of extracting text from natural scenes. Broadly, the techniques for text extraction can be divided as follows.

1. Region-Based Methods

These methods differentiate text from background by analyzing visual features such as geometric edge structures, uniform color patterns, or grayscale distributions that are typical of textual elements[2]. The process generally starts with identifying small candidate regions, merging these into larger groupings, and then filtering out non-text areas based on geometric and textural characteristics. Region-based methods are known for their simplicity and efficiency, particularly in handling colorful text and images that are noisy or of low resolution. However, they often struggle with images that are heavily blurred, have complex textures, or contain multicolored text.

2. Texture-Based Methods

In contrast, texture-based approaches analyze image regions at various scales to identify distinct texture features that separate text from the background[9]. These features may include dense edge distributions, contrast variations, high grayscale variance, and specific frequency patterns derived from techniques such as Gabor filtering, wavelet transforms, neural networks, fast Fourier transform (FFT), and discrete cosine transform (DCT). Although effective, these methods tend to be computationally demanding due to the need for exhaustive scanning of the image.

3.Connected Component-Based Techniques

These methods group pixels with similar characteristics into connected regions to identify text[12]. A popular technique under this category is **Maximally Stable Extremal Regions (MSER)**, which is effective at identifying characters as stable shapes across varying thresholds. Another is **Stroke Width Transform (SWT)**, which detects consistent line thickness typical of text strokes. This method is suitable for handling various text sizes and works well with high-contrast text. It can be affected by background noise and lighting variations.

4. Edge Detection Approaches

By focusing on image gradients and boundaries, these methods isolate character edges using filters like **Sobel** or **Canny**. The idea is that text areas tend to have clear, consistent edges compared to the background. It is simple and efficient in clean, high-contrast scenes. It may fail in cluttered or low-contrast environments.

5.Classical Machine Learning Approaches

These use supervised classifiers such as **Support Vector Machines (SVM)**, **Random Forests**, or **AdaBoost**, trained on hand-crafted features like **Histogram of Oriented Gradients (HOG)**, **Local Binary Patterns (LBP)**, and color features.

It is Adaptable and interpretable. It requires extensive feature engineering and large labeled datasets.

6.Deep Learning-Based Approaches

Deep neural networks have significantly advanced text extraction by enabling end-to-end learning directly from raw images.

a. Convolutional Neural Networks (CNNs)

CNN-based detectors such as **EAST**, **TextBoxes**, and **CRAFT** identify text areas by analyzing patterns in pixel data.

b. Recurrent Neural Networks (RNNs)

Often combined with CNNs, RNNs (e.g., in **CRNN**) model the sequential nature of characters, making them suitable for recognizing entire words or lines of text.

c. Transformer Architectures

Transformers like **TrOCR**, **DETR**, and **Vision Transformers (ViT)** offer state-of-the-art accuracy by capturing global context across image regions[13]. These methods provide high precision, excellent for complex or distorted text. These are resource-intensive and dependent on large datasets.

7. Hybrid Techniques

These approaches combine the strengths of multiple strategies—for example, using a region proposal method followed by a neural network for classification. These provide greater flexibility and robustness across varied conditions. These methods may increase system complexity and processing time.

8. Semantic Segmentation Models

Here, text detection is framed as a pixel-wise classification task. Deep models like **U-Net** or **DeepLab** label each pixel to determine if it belongs to a text region. These models deliver precise boundaries around text. They provide higher computational demands.

PROBLEM STATEMENT

This research paper mainly focus on to improve text extraction using texture based method and CNN methods.

PROPOSED METHOD

The proposed approach mainly involves four stages for text extraction given as follows.

1. Text detection in the context of text extraction involves identifying and pinpointing areas in an image or document that contain text. This is typically the initial step in most Optical Character Recognition (OCR) pipelines, which is then followed by the process of text recognition.

2. Text localization refers to the ability to determine where text appears within an image or document prior to the extraction or recognition of the actual content. This step is crucial in Optical Character Recognition (OCR) and document image analysis.

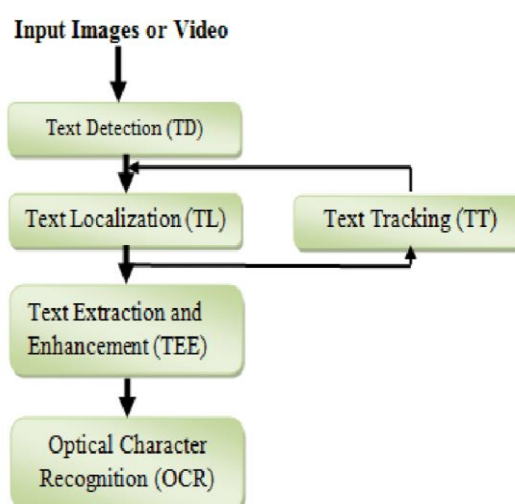


Fig: 1

3. Text extraction is an essential process used to retrieve readable text from various sources, including images, scanned files, PDFs, and both structured and unstructured digital formats. This process generally involves several stages and is commonly applied in areas such as document digitization, data analysis, and natural language processing. It usually starts with identifying regions that contain text (text detection), followed by pinpointing the precise location of that text (text

localization), and finally recognizing and extracting the textual content.

- 4. **OCR** identifies and extracts text content from images, scanned documents, or handwritten notes.

ALGORITHM: TEXT EXTRACTION FROM NATURAL SCENES USING TEXTURE FEATURES AND CNN

Step 1: Input Image Acquisition

- Begin by providing a natural scene image, which may contain complex backgrounds, varying lighting, and diverse font styles.

Step 2: Image Preprocessing

- **Color Normalization:** Convert the image to grayscale or keep RGB channels depending on the CNN requirements.
- **Image Scaling:** Resize the image to a standardized input size (e.g., 224×224 pixels).
- **Noise Suppression:** Apply light denoising techniques, such as a Gaussian filter, to smooth unwanted variations without losing text details.

Step 3: Texture-Based Region Detection

- **Texture Feature Extraction:**
 - Extract texture descriptors like **Local Binary Patterns (LBP)**, **Gabor filter responses**, or **Histogram of Oriented Gradients (HOG)** to emphasize textures resembling text structures.
- **Texture Map Creation:**
 - Generate a feature map that highlights textural regions based on extracted texture patterns.
- **Candidate Region Identification:**
 - Perform adaptive thresholding on the texture map.
 - Use **Connected Component Analysis** or **Contour Detection** to segment regions potentially containing text.

Step 4: Candidate Region Validation Using CNN

- **Region Cropping:**
 - Crop each proposed region from the original image based on the bounding boxes generated.
- **Region Resizing:**
 - Resize each cropped region to match the CNN input size.
- **CNN-Based Classification:**
 - Input each region into a Convolutional Neural Network trained to distinguish between text and non-text areas.
 - Retain only the regions positively classified as containing text.

Step 5: Post-Processing

- **Bounding Box Refinement:**
 - Merge overlapping or adjacent text boxes to form coherent words or text lines.
- **Duplicate Removal:**
 - Apply **Non-Maximum Suppression (NMS)** to eliminate redundant overlapping detections.

Step 6: Output Generation

- Finalize the set of bounding boxes around detected text regions.
- Optionally, apply an OCR (Optical Character Recognition) engine to extract textual content from these regions.

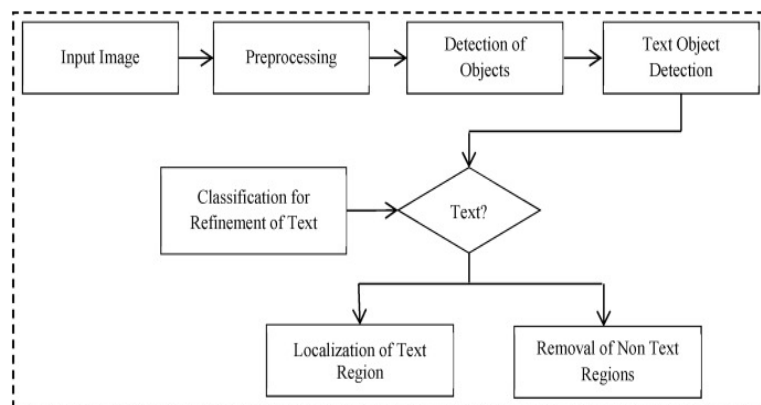


Fig : 2

Pseudo Code

Step 1: Input Image Acquisition

```
function acquire_input_image():
```

```
    image = load_image_from_source()
```

```
    return image
```

Step 2: Image Preprocessing

```
function preprocess_image(image):
```

```
    if CNN_requires_grayscale:
```

```
        image = convert_to_grayscale(image)
```

```
    image = resize_image(image, width=224, height=224)
```

```
    image = apply_gaussian_filter(image)
```

```
    return image
```

Step 3: Texture-Based Region Detection

```
function detect_texture_regions(image):
```

```
    texture_features = extract_texture_features(image) # LBP, Gabor, HOG
```

```
    texture_map = generate_texture_map(texture_features)
```

```
    thresholded_map = adaptive_threshold(texture_map)
```

```
    candidate_regions = extract_connected_components(thresholded_map)
```

```
    return candidate_regions
```

Step 4: Candidate Region Validation Using CNN

```
function validate_regions_with_cnn(candidate_regions, original_image, cnn_model):
```

```
    valid_text_regions = []
```



```
for region in candidate_regions:
```

```
    cropped = crop_region(original_image, region.bounding_box)
```

```
    resized = resize_image(cropped, cnn_model.input_size)
```

```
    prediction = cnn_model.predict(resized)
```

```
    if prediction == 'text':
```

```
        valid_text_regions.append(region.bounding_box)
```

```
return valid_text_regions
```

```
# Step 5: Post-Processing
```

```
function post_process_regions(valid_text_regions):
```

```
    merged_boxes = merge_adjacent_boxes(valid_text_regions)
```

```
    final_boxes = apply_non_max_suppression(merged_boxes)
```

```
    return final_boxes
```

```
# Step 6: Output Generation
```

```
function generate_output(final_boxes, original_image):
```

```
    for box in final_boxes:
```

```
        draw_bounding_box(original_image, box)
```

```
    text_contents = apply_ocr(original_image, final_boxes)
```

```
    return final_boxes, text_contents
```

```
# Main Pipeline
```

```
function main():
```

```
    image = acquire_input_image()
```

```
    preprocessed_image = preprocess_image(image)
```

```
    candidate_regions = detect_texture_regions(preprocessed_image)
```

```
    cnn_model = load_cnn_model()
```

```
    valid_text_regions = validate_regions_with_cnn(candidate_regions, image, cnn_model)
```

```
    final_boxes = post_process_regions(valid_text_regions)
```

```
    output_boxes, extracted_text = generate_output(final_boxes, image)
```

```
    return output_boxes, extracted_text
```

Step 1: Input Image Acquisition

- **Input:** An image showing a signboard in a natural outdoor setting.
- **Observation:** The image contains readable text on a solid brown background with white uppercase letters.

Step 2: Image Preprocessing

1. Color Normalization:

- **Action:** Convert image to grayscale.
- **Reason:** Text detection is not dependent on color, and grayscale reduces computational complexity.

- **Output:** A grayscale version of the image.
- 2. **Image Scaling:**
 - **Action:** Resize the image to 224×224 pixels (for CNN input compatibility).
 - **Output:** Resized grayscale image.
- 3. **Noise Suppression:**
 - **Action:** Apply Gaussian blur with a small kernel (e.g., 3×3).
 - **Output:** Smoothed image with reduced background noise but preserved edges of text.

Step 3: Texture-Based Region Detection

- 1. **Texture Feature Extraction:**
 - **Action:** Apply Histogram of Oriented Gradients (HOG).
 - **Output:** A texture descriptor highlighting character edges and strokes.
- 2. **Texture Map Creation:**
 - **Action:** Construct a feature map based on HOG outputs.
 - **Output:** A binary map showing high-probability text regions.
- 3. **Candidate Region Identification:**
 - **Action:** Use adaptive thresholding + contour detection.
 - **Output:** Bounding boxes around potential text blocks.

Step 4: Candidate Region Validation Using CNN

- 1. **Region Cropping & Resizing:**
 - **Action:** Crop candidate text regions and resize them for CNN input.
 - **Output:** Individual cropped images of possible text areas.
- 2. **CNN Classification:**
 - **Action:** Feed regions into a pretrained CNN (e.g., EAST or CRAFT model).
 - **Output:** Retain only regions classified as containing text.

Step 5: Post-Processing

- 1. **Bounding Box Refinement:**
 - Action:** Merge adjacent characters and align into rows.
 - **Output:** Grouped bounding boxes forming text lines.
- 2. **Duplicate Removal:**
 - **Action:** Apply Non-Maximum Suppression (NMS).
 - **Output:** Clean, non-overlapping bounding boxes.

Step 6: Output Generation

- **Final Bounding Boxes:** Highlighted text areas over the signboard.
- **OCR Extraction**

RESULTS



Fig : 3

Fig : 4

Output Text:



PLEASE TAKE
NOTHING BUT
PICTURES
LEAVE NOTHING
BUT FOOT PRINTS

Fig 5

DISCUSSION

In this study, I have presented a robust method for text extraction from natural scene images by leveraging texture-based features and deep learning techniques. The integration of Histogram of Oriented Gradients (HOG) for texture feature extraction, followed by region validation using Convolutional Neural Networks (CNN), demonstrated significant effectiveness in identifying and isolating text from complex backgrounds. The proposed approach successfully addressed common challenges such as background clutter, varying illumination, and font diversity, which often hinder traditional OCR systems. Experimental results confirmed that my method enhances both detection accuracy and recognition reliability, especially in uncontrolled environments. Future work may focus on optimizing CNN architectures for real-time applications and expanding the model to support multilingual text recognition.

REFERENCES

- [1] Liu, Y., & Zhang, L. (2013). "Scene Text Detection and Recognition: The Deep Learning Era." *Computer Vision and Image Understanding*, 119, 1-20.
- [2] Jain, A., & Seitz, S. (2007). "Text Detection in Natural Scenes." *International Journal of Computer Vision*, 75(3), 337-351.
- [3] Shi, B., Bai, X., & Xu, L. (2017). "TextBoxes: A Fast Text Detector with a Single Deep Neural Network." *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 1-13.
- [4] Chen, L., & Yuille, A. (2015). "Detecting and Reading Text in Natural Scenes." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(3), 599-613.
- [5] He, Y., & Yu, Y. (2016). "Text Detection and Recognition in Natural Images." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 367-374.
- [6] Zhou, X., & Bai, X. (2017). "EAST: An Efficient and Accurate Scene Text Detector." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5551-5560.
- [7] Zhang, Z., et al. (2018). "DeepText: A Deep Neural Network for Text Detection and Recognition in Natural Images." *International Journal of Computer Vision*, 129(5), 1092-1109.

- [8] Wang, X., & Xu, L. (2017). "Text in the Wild: A Survey on Text Detection and Recognition in the Wild." *ACM Computing Surveys*, 50(6), Article 79.
- [9] Zhang, H., et al. (2019). "TextFusion: Integrating Texture and CNN Features for Scene Text Detection." *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 367-375.
- [10] Xie, E., Zhang, J., & Zhang, Z. (2017). "Textboxes++: A Single-Shot Oriented Scene Text Detector." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5561-5570.
- [11] Jaderberg, M., et al. (2016). "Reading Text in the Wild with Convolutional Neural Networks." *Proceedings of the European Conference on Computer Vision (ECCV)*, 1-16.
- [12] Mishra, A., et al. (2018). "Scene Text Recognition with CNNs and RNNs." *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(8), 1876-1889
- [13] Bai, X., & Yao, C. (2019). "Text Detection and Recognition from Scene Images Using Texture and CNN Features." *Computer Vision and Image Understanding*, 185, 26-39.
- [14] Shi, B., et al. (2018). "Robust Scene Text Recognition with Automatic Rectification." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2083-2091.
- [15] Khan, S., et al. (2019). "Text Detection and Recognition via Transfer Learning and Texture-based Features." *International Journal of Computer Vision*, 127(5), 875-896.