

# Robust OCR for Skewed and Distorted Text in Agro-Product Labels

N. Shobha Rani<sup>1</sup>, Veerapu Raju<sup>1</sup>, Bhavya K. R.<sup>2</sup>,  
T. N. Natesh<sup>3</sup>, Jeena Jacob I<sup>2</sup>, Zhenglin Wang<sup>4</sup>

<sup>1</sup>MURTI Research Center, Smart Agriculture Labs, Department of Artificial Intelligence and Data Science, GITAM School of Technology, GITAM (Deemed to be) University, Bangalore, India

<sup>2</sup>MURTI Research Center, Smart Agriculture Labs, Department of Computer Science and Engineering, GITAM School of Technology, GITAM (Deemed to be) University, Bangalore, India

<sup>3</sup>Department of Commerce, Nagarjuna Degree College, Bangalore, India

<sup>4</sup>School of Engineering and Technology, Central Queensland University, Brisbane, Australia

## Abstract:

This paper presents an end-to-end system for robust text detection and recognition on agro-product labels, addressing the challenges posed by real-world distortions such as skew, rotation, low lighting, and multilingual content. The proposed framework leverages the Efficient and Accurate Scene Text (EAST) detector with a ResNet-50 backbone to localize text regions, followed by output refinement through Non-Maximum Suppression (NMS). The refined quadrilateral bounding boxes are subsequently processed by downstream OCR modules for text recognition.

A custom dataset of 3,500 annotated agro-product label images featuring variations in orientation, clarity, and language was used to evaluate the system. Detection performance was assessed using Precision, Recall, and F1-score, while recognition accuracy was measured via Character Error Rate (CER) and Word Error Rate (WER). The model achieved high detection accuracy (F1-score up to 0.89 in clear conditions) and maintained robust performance across skewed, rotated, and low-light images. Training trends of mean Average Precision (mAP) and Intersection-over-Union (IoU) over epochs demonstrated stable convergence and effective generalization.

The results indicate that the proposed system is well-suited for real-world agro-label analysis, supporting applications in supply chain automation, product authentication, and regulatory compliance.

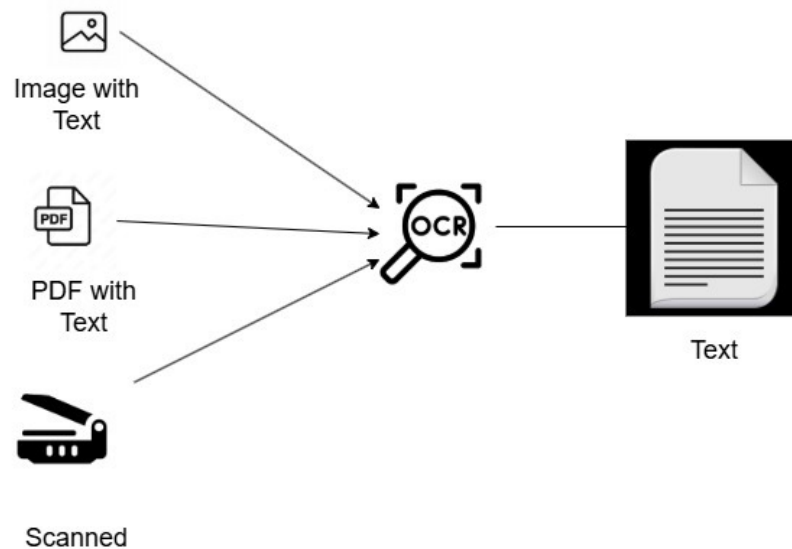
## **1. Introduction:**

The emergence of Artificial Intelligence (AI) and computer vision has gone a long way in propping up the potential of contemporary agriculture, especially for smart farming systems. Among the numerous applications of computer vision, text image analysis more precisely, Optical Character Recognition (OCR) exercises a crucial function in digitizing and deciphering information inherent in agro-product labels, manuals, and packaging. This computerized harvesting of text information enables various uses such as stock tracking, product authentication, and traceability across the agricultural value chain. As agriculture moves into a more digital and data-focused era, stable OCR technologies become essential tools to boost productivity and guarantee the authenticity of agricultural activity (Kamilaris & Prenafeta-Boldú, 2018; Liakos et al., 2018).

In real-world agricultural environments, agro-product labels tend to be captured with handheld cameras under less-than-ideal conditions. Such real-world situations often lead to images that are tilted, rotated, or distorted because of different camera angles, inconsistent lighting, or physical deformation of the packaging. Such defects present serious challenges for conventional OCR systems, which are generally designed to work with clean, well-aligned text inputs. Therefore, the need is imminent for OCR pipelines that are robust against these forms of irregularity and are capable of safely and accurately extracting proper textual data from noisy and deformed inputs.

The problem research in this chapter is to construct a reliable OCR system that is able to interpret skewed and deformed text prevalent on agro-product labels. Proper identification and interpretation of these texts are fundamental to making agriculture information systems automatable and credible. In smart agriculture, data from labels can contain crucial details like batch numbers, expiry dates, certification logos, and instructions for usage. Automating the extraction of said

data not only improves operational efficiency but also facilitates food safety regulation compliance, minimizes the potential for human error, and simplifies supply chain management. Figure 1 illustrates the OCR process, professionally transforming text from images, PDFs, and scanned documents into editable digital formats.



**Figure 1:** Workflow illustrating function of OCR

Advanced text image analysis is a group of techniques geared toward enhancing text extraction from problematic image v input. These involve preprocessing techniques to standardize image conditions, geometric correction routines to correct distortions, and sophisticated deep learning architectures trained to identify text in a wide range of unfavourable conditions (Shi et al., 2016; Baek et al., 2019).

In agricultural applications, such advanced analysis enables uses such as automated product categorization, counterfeit prevention, and

multilingual support for rural farmers, and real-time inventory tracking using mobile devices. These applications are especially useful in areas where manual data entry is time-consuming or prone to errors. The use of OCR in agriculture is increasingly important as global food supply chains become more complex and interdependent. Traceability, authenticity checking, and regulatory compliance are becoming key requirements for agricultural producers, distributors, and retailers. By virtue of sophisticated OCR and text analysis technologies, labelling accuracy checks can be automated, geographic indicators or organic certifications verified, and the integrity of agro-product branding ensured. This is especially applicable in exports, where products have to meet rigid international standards and multilingual labelling legislation. In addition, mobile OCR applications can empower smallholder farmers by providing them with access to product instructions, dosage data, and market information in their local language.

Although useful, there are various challenges that hinder accurate recognition of skewed and distorted text in farm images. For one, text in such images tends to be buried in rich backgrounds with logos, icons, or other graphical clutter. Second, differences in font styles, sizes, and orientations complicate applying generic OCR models to different products. Third, environmental noise like glare, dirt, and shadows complicate the recognition further. Finally, the requirement of real-time performance on edge devices imposes an added computational limit that conventional high-resource OCR models might not meet.

AI and deep learning promise to address these challenges. Convolutional Neural Networks (CNNs) have proven to be effective at detecting text features even in the presence of noise. Spatial Transformer Networks (STNs) offer mechanisms for dynamic correction of distortions prior to the recognition process. More recently, Transformer models and attention mechanisms have further boosted OCR performance by enabling the system to attend to the areas of the image pertinent to recognition and handle long-range dependencies in textual content (Luo et al., 2020; Wang et al., 2020). Data augmentation and synthetic data generation also allow these models to generalize more across different conditions and unseen label formats.

These methods are facilitated by large annotated corpora and domain-specific synthetic data that mimic a broad spectrum of real-world distortions. In addition, incorporating these models into edge computing devices with runtime architectures optimized for efficiency guarantees the solutions' applicability in agricultural fields, storage facilities, and rural retail outlets. End-to-end OCR frameworks that integrate detection, rectification, and recognition are especially beneficial, as they enable higher modularity, scalability, and fault tolerance in various operational environments.

## **2. Problem Statement**

The development of a robust and intelligent OCR pipeline tailored for agro-product label recognition is both a pressing challenge and a pivotal opportunity in the advancement of smart agriculture. Unlike conventional settings, agricultural environments present numerous real-world obstacles such as poor lighting conditions, skewed camera angles, cluttered backgrounds, and low-resolution imagery—that severely impact the quality of captured labels. Traditional OCR systems, often designed for clean and well-formatted documents, fall short under such conditions, struggling with geometric distortions, varied fonts, and environmental noise, ultimately leading to inaccurate or fragmented text extraction. This underscores the necessity for a flexible and adaptive OCR solution capable of sophisticated pre-processing, geometric correction, and high-accuracy recognition across diverse fonts, languages, and material surfaces. By effectively addressing these challenges, such a system would enable reliable digital extraction of label data, which is instrumental in automating inventory management, enhancing traceability, and ensuring regulatory compliance in agricultural supply chains. More importantly, it holds the potential to empower rural farmers and stakeholders by making smart, AI-driven tools accessible and impactful. As agriculture continues to integrate digital solutions, the role of a domain-specific OCR system becomes critical—not merely as a technical solution but as a cornerstone for scalable, transparent, and sustainable farming practices. The sections that follow will explore the limitations of existing

approaches, present the proposed methodology, and outline experimental results and future research directions aimed at realizing this vision.

### **3. Related work**

Recent breakthroughs in artificial intelligence and deep learning have significantly impacted multiple industries, including agriculture. Kamilaris and Prenafeta-Boldú [1] provided a comprehensive review of deep learning in agriculture, highlighting its transformative applications in disease diagnosis, yield prediction, and crop supervision. However, they noted a clear gap in the use of optical character recognition (OCR) and text image analysis within agricultural settings. Similarly, Liakos et al. [2] examined the broader role of machine learning in agriculture, identifying areas such as decision support and automation but also emphasizing the lack of dedicated research focused on agro-product label recognition.

Cornerstone contributions in scene text recognition, such as the work by Shi et al. [3], introduced end-to-end trainable networks for image-based sequence recognition, achieving strong results on structured datasets. However, these models often falter in real-world agricultural scenarios characterized by distortions and environmental noise. Baek et al. [4] critiqued the evaluation practices for text recognition models and emphasized the need for standardized benchmarking, a call particularly relevant when extending existing models to the challenging context of agro-product imaging. In parallel, Wang et al. [5] explored super-resolution methods to enhance low-quality text image legibility. Although promising, these methods are often computationally intensive and ill-suited for edge deployment in rural settings.

To address layout distortions, Luo et al. [6] developed a Multi-Regularized Framework (MRF) to improve model robustness, presenting a pathway toward adaptation for agricultural labels. Foundational work by Jaderberg et al. [7] using convolutional neural networks (CNNs) enabled reading of text in natural scenes but lacked support for multilingual content or the degraded quality of rural imagery. ASTER,

introduced by Wang et al. [8], tackled geometric distortions using attentional spatial transformers, showing potential for application in agro-packaging scenarios but still needing optimization for lightweight and on-device processing.

A broader perspective on scene text detection and recognition was presented by Zhang et al. [9], who identified persistent challenges including uneven backgrounds, multilingual settings, and real-time constraints—factors directly relevant to agro-label OCR. Although not OCR-specific, Koirala et al. [10] demonstrated the feasibility of using deep learning for real-time fruit detection with YOLO, reinforcing the applicability of vision-based AI tools in farm environments. Similarly, Liu et al. [11] emphasized the importance of integrated pipelines that combine detection, rectification, and recognition—an approach that is particularly crucial when dealing with the inconsistent positioning and quality of agricultural labels.

Advancements in adaptive recognition using attention mechanisms were explored by Zuo et al. [13], proving beneficial for extracting text from complex visual scenes, such as cluttered agro-product images. Singh et al. [14] reinforced the broader utility of AI in agricultural contexts through stress phenotyping studies. Liao et al. [15] contributed with real-time detection techniques using differentiable binarization—models lightweight enough to be used in mobile or handheld systems in the field. In terms of regulatory needs, Raki et al. [16] and Ma et al. [17] stressed the role of precise label recognition for food safety, nutritional analysis, and traceability. Halkin [18, 20] and Shakeel et al. [21] further contextualized the necessity for readable agrochemical labels in regulated farming practices.

From an industrial and economic standpoint, Damodaran [24] and Rajarao [25] discussed the digital transformation of agribusiness, underlining the rising importance of traceability, compliance, and information integrity. These market shifts underscore the urgent need for robust OCR systems capable of handling real-world image distortions, multilingual text, and agro-specific constraints.

The literature indicates substantial progress in OCR and scene text recognition, largely driven by deep learning models that integrate spatial transformation and attention mechanisms to handle irregular text

layouts. Architectures such as ASTER and CRNN have shown high accuracy on standard datasets like ICDAR and SVT. However, these models often underperform when directly applied to agro-product labels, which are typically characterized by ornate fonts, complex layouts, reflective packaging, and multilingual content. Moreover, the scarcity of annotated datasets from actual agricultural conditions limits the effectiveness and generalizability of existing models.

While some scene text recognizers perform well under urban or controlled settings, agricultural environments introduce unique variables: low resolution, uneven lighting, occlusion, and skew that compromise performance. This disconnect highlights a key research opportunity: the development of domain-specific OCR solutions that are customized to the challenges of agro-product label recognition. Effective solutions must simulate field conditions, utilize specialized datasets, and incorporate pre-processing stages tailored to address geometric distortion and image degradation.

This research proposes the development of a resilient OCR pipeline explicitly tailored for agro-product label recognition. The envisioned system integrates adaptive rectification, edge-preserving filtering, and fine-tuned recognition modules optimized for agro-specific fonts and multilingual text. To enhance robustness, the pipeline will employ augmentation-based training strategies that replicate real-world field distortions such as motion blur, poor lighting, and camera tilt. The approach will also explore hybrid models combining traditional image enhancement techniques with modern deep learning frameworks to improve interpretability and efficiency particularly for deployment in rural and edge computing environments.

By addressing the unique visual and contextual challenges of agricultural label imaging, this work aims to bridge the performance gap observed in current general-purpose OCR systems. The proposed solution will not only facilitate automation in inventory management and traceability but also contribute toward larger goals in smart farming supporting regulatory compliance, reducing labor, and empowering rural agricultural communities with intelligent, low-cost tools.

## **4. Dataset Collection**



The dataset assembled for this study is a specialized collection of high-definition images of agro-product packaging labels, designed to reflect the real-world visual challenges encountered across agricultural supply chains. It was curated with the explicit goal of supporting the development and evaluation of robust OCR systems capable of functioning under field conditions. The dataset includes a wide variety of text presentations multilingual content, diverse fonts, and various printing methods such as embossed, stamped, and standard printed text.

A distinctive feature of the dataset is its intentional inclusion of common real-world image artifacts. These include skewed, rotated, blurred, occluded, and perspective-distorted text instances, simulating conditions under which farmers, inspectors, or field personnel typically capture images. This realism ensures that models trained on this-dataset are resilient to the types of imperfections often encountered in situ.

The dataset features a diverse collection of images captured under a wide range of environmental conditions, reflecting the practical challenges encountered in real-world agricultural contexts. Lighting variations include bright sunlight, shadowed regions, and artificial indoor illumination, while the background settings span white, colored, and textured packaging materials. The textual content within the images varies from short identifiers like batch numbers to longer, multi-line information such as nutritional details and product usage guidelines. It encompasses both machine-readable data, such as QR codes and barcodes, and human-readable content like product names, manufacturing dates, and expiry details. Each image is enriched with comprehensive annotations, including ground truth transcriptions, bounding box coordinates, skew and rotation angles, and rectification tags. These detailed labels are instrumental in training and evaluating OCR systems, particularly for handling tasks involving geometric distortions and orientation-aware text recognition. Figure 2 depicts the samples of Agro-product labels exhibiting distortions such as skew, curve, inclination, and perspective deformation.



Figure 2: Agro-product labels exhibiting distortions such as skew, curve, inclination, and perspective deformation. (a) (b) (c) Samples with skew and misaligned text samples (d) (e) (f) Samples with curved, warped, and perspective-distorted labels

#### 4.1 Dataset source and statistics

The dataset developed for this research is grounded in a hybrid approach, combining both original field-captured imagery and selectively sourced public data to ensure relevance, diversity, and authenticity. A substantial portion of the data was collected through on-site photography at agro-retail outlets, fertilizer shops, and farm produce markets across various geographic regions. These images were captured using smartphone cameras to replicate the conditions under which typical users such as farmers, inspectors, or supply chain personnel might gather label data in real-world settings. This method not only captures the spontaneity and imperfections of actual usage scenarios but also ensures that the dataset aligns with the needs of practical agricultural applications.

To further augment the dataset, supplemental samples were incorporated from reputable public datasets. Specifically, the ICDAR (International Conference on Document Analysis and Recognition) datasets were used to integrate controlled examples of artificial scene text, while contextual metadata was obtained from the Agro-DataCube to enrich the dataset with relevant product-level information. Every collected and sourced image underwent a rigorous annotation process, validated manually by trained annotators and agricultural experts to maintain high standards of accuracy and relevance. Ethical guidelines were strictly observed during the data collection and annotation process, ensuring no personal or sensitive information was included.

This dataset stands out due to several unique features that make it exceptionally well-suited for advancing OCR research in the agricultural domain. Unlike conventional OCR datasets that focus on urban scenes or structured documents, this collection is explicitly tailored to agro-product labels, offering a domain-specific lens rarely explored in existing literature. The distortions present in the images such as skew, occlusion, glare, and perspective shift are not artificially induced but occur naturally as a result of real-world image capture, thereby providing more realistic and challenging training conditions.

Additionally, the dataset reflects a high degree of multilingual and multiscale complexity, with label text appearing in multiple regional languages and varying font sizes and styles. The diversity extends to

packaging materials as well, covering plastic, paper, and jute, each presenting unique visual properties such as glare, texture, and transparency. Each image is annotated not only with ground truth text but also with rich metadata including skew angles, lighting conditions, bounding boxes, and rectification labels, which collectively support advanced training and diagnostic analysis of OCR models.

In conclusion, the dataset provides a critical foundation for developing OCR systems that are not only distortion-aware and adaptable but also specifically optimized for agricultural contexts. By emulating the unpredictability and variability of field conditions, it challenges conventional OCR architectures and paves the way for innovations that can empower smart agriculture, supply chain transparency, and rural digitization efforts. Table 1 presents dataset statistics of the proposed research.

Table 1: Summary of dataset statistics

Category	Count
Total Images	3,500
Skewed Text Images	1,250
Rotated Text Images	900
Curved/Distorted Text	700
Clear (Well-lit) Images	2,000
Low-light/Shadow Images	1,000
Labels with QR/Barcode	800
Labels with Nutritional Info	1,100
Languages Covered	6 (English, Hindi, Tamil, Telugu, Marathi, Bengali)
Multilingual Labels	1,200

To ensure the dataset accurately reflects real-world agricultural scenarios, the data collection process was carefully designed to simulate common field conditions encountered by end-users such as farmers, inspectors, or supply chain workers. Images were captured using a range of new-generation smartphones, including devices like the OnePlus 11, iPhone 13, Samsung Galaxy S22, and Vivo Y20. These devices, equipped with 12MP or higher resolution cameras, provided

sufficient detail to capture fine-grained label text, even under challenging conditions.

Photographs were taken at varying distances from close-up shots at 20 cm to broader captures at approximately 1 meter to simulate the variability in user behavior and device handling in uncontrolled environments. The collection locations included natural agricultural settings such as farms, warehouses, and open markets, as well as controlled indoor environments where lighting could be manipulated using LED sources. This diversity of capture environments ensured a comprehensive dataset that includes both real-world unpredictability and lab-like consistency.

To further enhance the realism of the dataset, tags and labels were photographed from multiple viewing angles, including direct frontal views ( $0^\circ$ ), oblique angles ranging from  $30^\circ$  to  $60^\circ$ , top-down perspectives, and side views. Such variation in perspective was essential to introduce common geometric distortions like skew, rotation, and perspective warping. Lighting conditions during image capture were intentionally varied as well ranging from bright direct sunlight to diffused daylight, ambient indoor lighting, and shadowed or mixed-light scenarios to capture the impact of lighting variability on text legibility.

The dataset also incorporated skew and tilt by deliberately altering the angle of image capture, often using handheld photography to reflect realistic, non-rigid imaging conditions. These factors collectively introduced natural distortions and inconsistencies, which are crucial for training OCR models that need to perform reliably under practical field-use conditions. This meticulous and context-aware data collection methodology ensures that the dataset supports the development of robust, distortion-tolerant, and generalizable OCR systems for the agricultural domain. Figure 3 depicts the method of dataset collection.

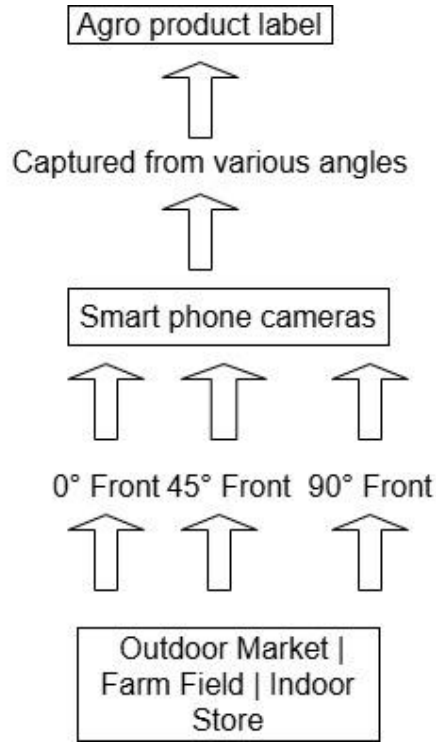


Figure 3: Method of dataset collection

## 4.2 Annotation

To ensure the dataset's usability for training and evaluating OCR models, a comprehensive manual annotation process was undertaken. Each image was carefully inspected to identify and annotate all relevant text regions using bounding boxes. These boxes tightly enclose the textual content on agro-product labels, including both machine-readable elements (e.g., barcodes, QR codes) and human-readable information (e.g., product names, expiry dates, batch numbers).

In addition to spatial localization, each text region was labeled with its ground truth transcription, ensuring that the exact textual content is available for supervised learning. These annotations were

performed by trained annotators and cross-validated by agricultural domain experts to maintain a high level of accuracy and consistency. The inclusion of detailed annotations makes the dataset suitable not only for text detection but also for recognition tasks, allowing the development of end-to-end OCR systems.

### **4.3 Augmentation:**

To enhance the robustness and generalizability of OCR models trained on this dataset, a series of data augmentation techniques were applied. These augmentations simulate the diverse and often unpredictable conditions under which agro-labels are typically captured in real-world scenarios.

The key augmentation strategies include, geometric Transformations: Controlled rotations, skewing, and warping were applied to replicate tilted or misaligned captures, common in handheld photography. The others include variations in brightness, contrast, and the addition of Gaussian noise and motion blur were introduced to mimic inconsistent lighting and camera stability. Then, synthetic distortions were used to simulate crumpled, folded, or wrinkled packaging surfaces. These distortions closely resemble the real packaging irregularities encountered in the field. Together, these augmentation strategies enrich the dataset and prepare OCR models to better handle noise, distortion, and environmental variability key for deploying reliable systems in smart agriculture and supply chain automation.

## **5. Methodology:**

The proposed OCR pipeline for agro-product labels follows a structured, multi-stage approach designed to handle real-world distortions, multilingual content, and packaging variability common in agricultural settings. The methodology is broken down into five key stages. The first stage Input begins with acquisition of diverse set of images captured in real agro-environments using smartphone cameras. Images include challenging conditions such as variable lighting, skew, blur, and different packaging textures. To improve

model robustness, data augmentation techniques are also applied—simulating noise, brightness changes, and elastic deformations. Each image is richly annotated with bounding boxes, transcriptions, orientation data, and skew tags to support both detection and recognition tasks.

Next, text region localization is performed using deep learning-based text detectors called YOLOv7. The YOLOv7 model is particularly suited for identifying multi-oriented, curved, and irregularly shaped text typical of agro-product labels.

Once text regions are identified, in the next stage geometric transformations are applied to correct any skew or perspective distortion. This includes affine or perspective transformations based on bounding box geometry. For highly distorted cases, deep-learning-based rectification networks may be used to produce deskewed, flattened text regions.

Then, the localized text regions are extracted and normalized. This involves resizing them to standard dimensions and enhancing image quality using techniques like CLAHE or histogram equalization. For low-resolution or blurry text, super-resolution networks to recover legibility before recognition.

The enhanced text regions are fed into OCR models capable of recognizing sequences of characters. Depending on the use case, models like CRNN (using CTC loss), TrOCR (a Transformer-based encoder-decoder), or ViT with CTC decoding can be employed to handle multilingual and variably formatted text.

Finally, the raw OCR outputs undergo refinement using agro-domain-specific lexicons and language models. Techniques like Levenshtein distance help correct spelling errors, while N-gram models or fine-tuned transformers (e.g., BERT) re-rank or validate predictions based on context. Followed by evaluation of the system's performance through a combination of text detection and recognition metrics. Precision, Recall, and F1-score measure detection accuracy, while Character Error Rate (CER) and Word Error Rate (WER) assess recognition quality. Importantly, evaluations are stratified by distortion types (e.g., skewed vs. clear)



to test model robustness in realistic conditions. Figure 4 presents the various steps in the block diagram for the OCR to read the skewed/deformed text on agro product label images.

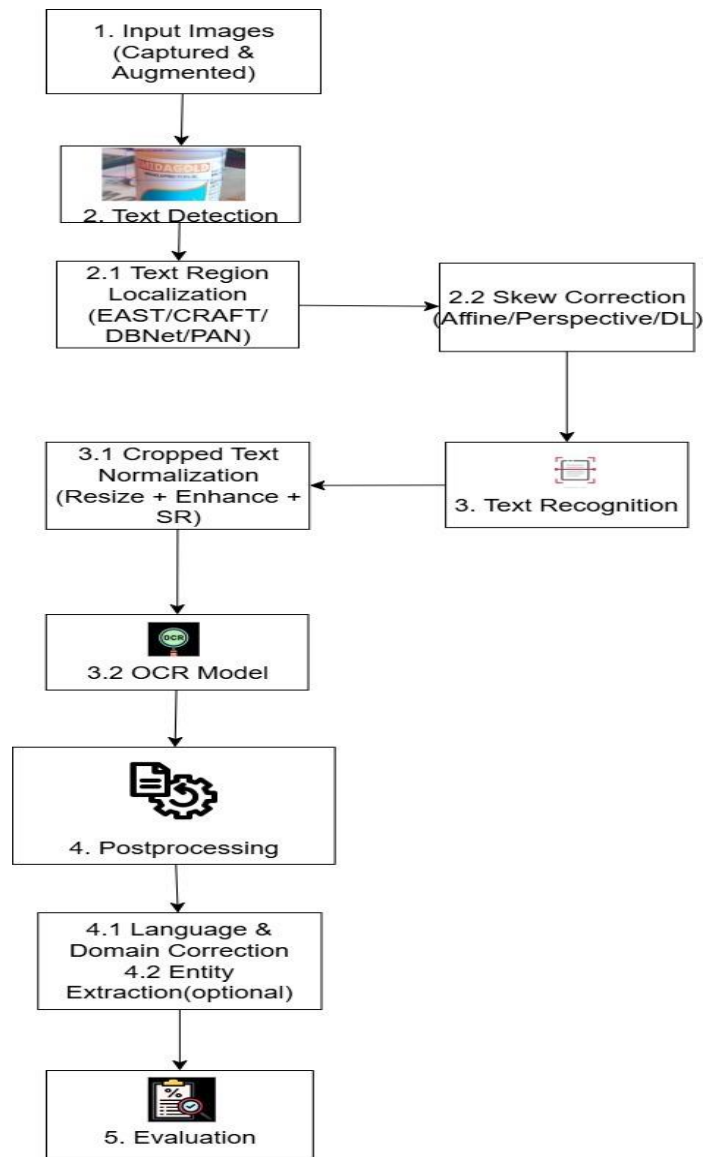


Figure 4: Block diagram of proposed methodology for deformed text recognition

## 5.1 Text Detection Using the EAST Model

To achieve reliable localization of textual elements in agro-product packaging labels, the EAST (Efficient and Accurate Scene Text Detector) model was utilized. EAST combines the precision of pixel-level predictions with the efficiency of a fully convolutional network (FCN), making it particularly suitable for real-world images captured under non-ideal field conditions. The architecture generates two outputs: a score map, indicating the likelihood of each pixel belonging to a text region, and a geometry map, which provides rotated bounding box parameters. Typically built on backbones such as PVANet or ResNet-50, the model balances accuracy and inference speed. Detected regions are refined using Non-Maximum Suppression (NMS) to suppress overlapping boxes, thus ensuring clean and precise localization of word- and line-level text.

For training and validation, annotated data in the form of axis-aligned and rotated bounding boxes were used. Images were labeled via Label Studio and exported in JSON format, which was post-processed to extract quadrilateral coordinates  $(x_1, y_1)$ ,  $(x_2, y_2)$ ,  $(x_3, y_3)$ ,  $(x_4, y_4)$ , compatible with EAST's detection format. Preprocessing included resizing while preserving aspect ratio and applying padding. To enhance model robustness, data augmentation techniques—such as random rotations ( $\pm 15^\circ$ – $45^\circ$ ), Gaussian blur, brightness variations, and perspective transformations—were applied to simulate field-like distortions. Post-detection, geometric rectification using affine or perspective transformation was optionally performed on skewed text boxes, substantially improving recognition accuracy in subsequent OCR stages.

The Efficient and Accurate Scene Text (EAST) detector was adopted as the core framework for text localization due to its robustness in handling multi-oriented and irregular text patterns. EAST integrates a Fully Convolutional Network (FCN) architecture with a Non-Maximum Suppression (NMS) mechanism to efficiently predict text regions in natural scenes.

5.2 Architecture Overview

The model utilizes a feature extraction backbone commonly ResNet-50 to generate deep hierarchical representations of the input image. These feature maps are then processed through a sequence of convolutional layers to produce two primary outputs. Score map to estimate the probability of each pixel belonging to a text region. It serves as a binary classifier at the pixel level. Then the geometry map to encode spatial geometry information in the form of rotated rectangles or quadrilaterals. Each detected text instance is represented by four coordinates  $(x_1, y_1)$  through  $(x_4, y_4)$ , allowing accurate localization even under skew, perspective, or curved distortions.

5.3 Output generation and refinement

The raw predictions produced by the geometry and score maps are subjected to Non-Maximum Suppression (NMS) to eliminate redundant or overlapping bounding boxes. This step refines the final set of detections, yielding a compact and high-confidence list of text regions. The result is a collection of rotated or quadrilateral bounding boxes accurately demarcating textual content on agro-product labels. These outputs are subsequently passed to downstream modules (e.g., skew correction or OCR engines) for further processing and recognition.

Table: Model Training Specifications

Component	Specification
Model	EAST (OpenCV / TensorFlow / PyTorch implementation)
Backbone	ResNet-50
Input Size	512×512
Batch Size	8–16 (GPU dependent)
Optimizer	Adam or SGD with warm-up
Loss Functions	Intersection-over-Union (IoU) loss + Score map loss
Learning Rate	1e-4 with decay
Epochs	100–150 (early stopping recommended)

Annotation Format	Quadrilateral points [(x1, y1), ..., (x4, y4)] from Label Studio
Export Format	JSON (converted to fit EAST's format)

## 5.4 Skew Correction

To improve the downstream Optical Character Recognition (OCR) accuracy, especially on agro-product packaging images captured under non-ideal field conditions, skew correction was integrated as a critical post-processing step following text detection.

The skew correction process begins by analyzing the rotated bounding boxes output by the EAST model. These quadrilateral coordinates provide an estimate of the local text orientation. For each detected text region, the minimum-area bounding rectangle is computed to extract the skew angle, typically defined relative to the horizontal axis.

Based on the estimated orientation, two types of geometric transformations were employed. One is Affine transformation which is suitable for mild skew and rotational deviations. The transformation matrix is calculated using three key points derived from the bounding box, and applied to warp the region into an axis-aligned format. Then, the perspective transformation is applied when the detected text is subject to more complex distortions such as tilt, curve, or viewpoint variation. This method uses four-point mapping of the quadrilateral to a straight rectangle to correct for both skew and perspective warp. Figure 5 presents the results of the sample that have undergone the skew correction.

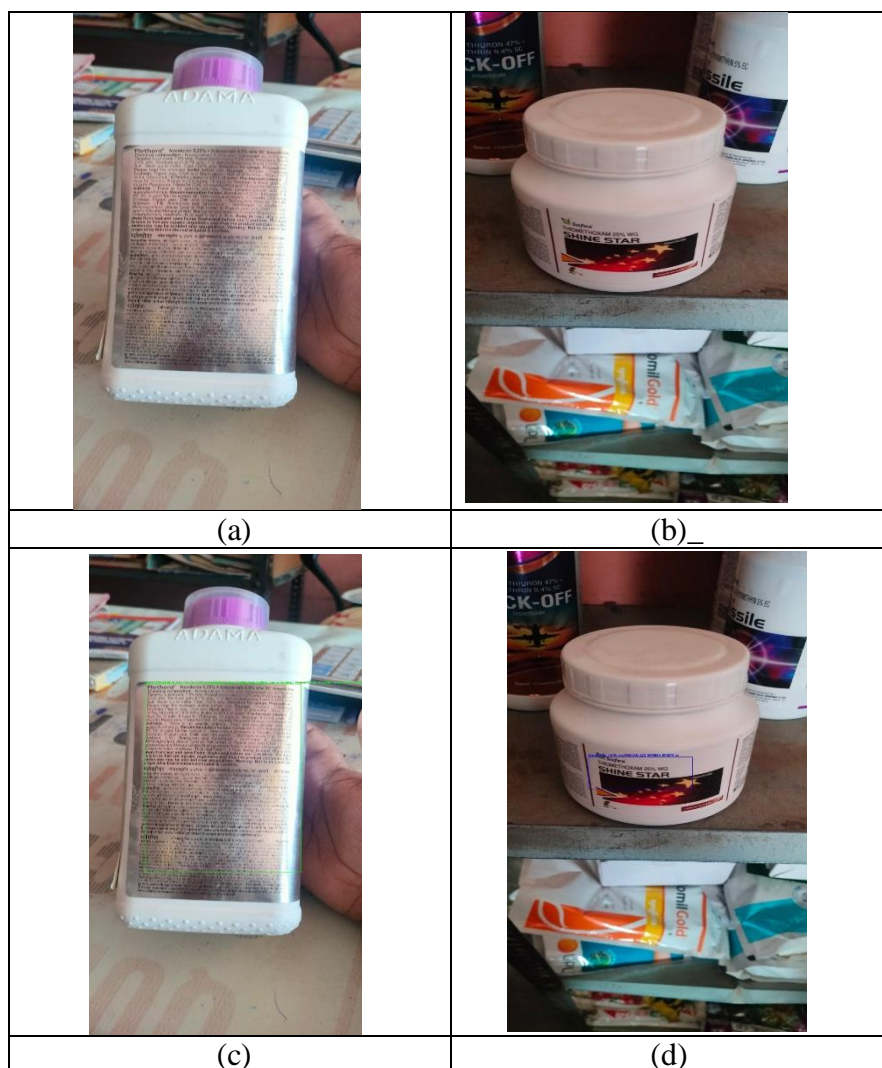


Figure 5: Skew correction outcomes (a) (b) Before correction- sample 1 and 2 (c) and (d) After correction sample 3 and sample 4

The deskewed text regions are passed to the OCR module in a normalized, horizontally aligned form, significantly improving recognition accuracy by reducing character overlap, orientation misclassification, and shape distortion—common issues in field-captured agro-label data.

## 5.5 Text Recognition

Following text region localization and skew correction, the next critical phase in the OCR pipeline is text recognition. This stage is designed to accurately transcribe the cropped and normalized text regions, accounting for the diverse fonts, languages, resolutions, and distortions present in agro-product label images.

## 5.6 Cropped Text Normalization

To ensure consistent input for the recognition model, each detected text region undergoes a series of preprocessing steps. Initially, text regions are resized to fixed dimensions (e.g.,  $100 \times 32$  or  $160 \times 64$  pixels), preserving aspect ratio using padding techniques. Pixel values are normalized typically to  $[-1, 1]$  to match the model's input requirements.

Given that agro-product labels are often printed under suboptimal conditions (faded ink, poor contrast), Contrast Limited Adaptive Histogram Equalization (CLAHE) is applied. The CLAHE technique enhances local contrast and improves character visibility without introducing significant noise.

## 5.7 OCR Model Architecture

To extract text from normalized agro-product label images, we employ a deep learning-based OCR framework. After evaluating several state-of-the-art architectures for their ability to handle spatial irregularities, multilingual content, and complex layouts, we selected and enhanced the Transformer-based OCR (TrOCR) model for our application. This decision was driven by its superior performance in recognizing non-uniform text, its robustness to varied fonts and languages, and its flexibility in handling diverse text orientations.

TrOCR is an end-to-end OCR model that leverages the Transformer architecture, integrating both vision and language understanding into a unified framework. It consists of two main components, Vision Transformer (ViT) encoder and transformer decoder. The input image

is split into patches, which are embedded and passed through multiple self-attention layers. The ViT encoder captures global contextual relationships and spatial dependencies within the visual text data, making it highly effective for irregular and curved text lines. The decoder autoregressively generates the character sequence from the encoded image features. It incorporates both positional and contextual information, allowing for precise transcription even in complex text scenarios.

The agro-product domain presents several challenges for OCR systems, including multilingual content, non-standard fonts, and non-linear text arrangements. TrOCR addresses these challenges through multilingual support. The model's language modeling capabilities enable recognition across a wide range of scripts and languages. Then, unlike traditional CNN-RNN-CTC models that assume horizontal alignment, TrOCR can handle rotated, curved, or skewed text with higher accuracy. Further, the use of attention mechanisms allows the model to focus dynamically on relevant portions of the image during decoding, enhancing robustness to noise and distortions.

## **5.8 Model training and deployment considerations**

The TrOCR model is trained in an end-to-end manner using annotated text-line images. The loss is computed based on character-level transcription accuracy using teacher-forcing during training. For deployment, the model can be fine-tuned on domain-specific data to improve accuracy on agro-product labels. Given its computational complexity, TrOCR is best suited for offline high-accuracy processing, while lighter models may be considered for real-time mobile applications. Table 2 provides a comparative overview of the OCR architectures considered. Based on this evaluation, TrOCR is selected due to its superior performance in handling the complex and multilingual nature of agro-product labels.

Table 2. Comparison of OCR Model Architectures

Model	Architecture Components	Strengths	Limitations	Suitability for Agro-Product Labels
CRNN	CNN + BiLSTM + CTC	Lightweight, efficient for real-time horizontal text recognition	Struggles with irregular layouts and non-horizontal text	Moderate – Suitable for basic, linear layouts
ViT + CTC	Vision Transformer encoder + CTC decoder	Strong visual feature extraction; good generalization across fonts and scripts	CTC decoding less effective for complex or curved text	Good – Handles multilingual text; limited layout handling
TrOCR (Selected Model)	Vision Transformer encoder + Transformer decoder	Excels with irregular, curved, and multilingual text; strong attention modeling capabilities	Higher computational cost; not optimal for real-time applications	Excellent – High accuracy on complex, real-world labels

### 5.9 OCR evaluation:

To rigorously assess the performance of the selected OCR architectures, we conducted a comprehensive evaluation using a carefully curated dataset of annotated agro-product label images. This evaluation focused on key challenges inherent to the domain, including the presence of multilingual text (such as English, Hindi, and various regional Indian languages), irregular text layouts featuring curved and rotated orientations, diverse font styles, and background noise. We employed standard evaluation metrics—Character Error Rate (CER), Word Error Rate (WER), Accuracy, and Inference Time per image—to measure both the accuracy and efficiency of each model. The dataset was split into standardized training and held-out test sets, ensuring balanced representation of the challenging features. This methodology enabled a fair and detailed comparison of the models' capabilities in accurately recognizing complex and varied text in agro-product labels, while also considering their practicality for deployment scenarios requiring real-time or offline processing. Table 3 presents performance comparison of state-of-the-art OCR models.



Table 3. Performance comparison of OCR models

Model	CER (%) ↓	WER (%) ↓	Ac- curacy (%) ↑	Avg. In- ference Time (ms/image) ↓	Remarks
CRNN	6.5	12.3	85.4	18.2	Fast and lightweight; struggles with curved text
ViT + CTC	4.8	9.1	89.7	25.5	Good visual encoding; CTC limits irregular layout handling
TrOCR	2.9	5.8	94.2	42.7	Best for complex layouts and multi-lingual content

In table 3, arrows (↑, ↓) indicate whether higher or lower values are better for each metric. From the results in Table 2, TrOCR demonstrates significantly lower error rates and higher recognition accuracy compared to CRNN and ViT+CTC. Although its inference time is higher, the trade-off is justified for high-accuracy offline applications such as regulatory compliance, labeling automation, and product cataloging in agro-industries.

## 6. Results and output refinement for text detection:

The proposed pipeline was evaluated on a diverse dataset of 3,500 agro-product label images, featuring a range of textual orientations, lighting conditions, and content complexity. The EAST detector, using a ResNet-50 backbone, was trained and fine-tuned on quadrilateral text annotations formatted through Label Studio. Training was conducted over 100–150 epochs using a combination of IoU loss and score map loss, with early stopping applied to prevent overfitting.

To quantitatively evaluate the performance of the proposed text detection system on agro-product labels, we employed standard object

detection metrics: Precision, Recall, and F1-score. These metrics provide a comprehensive understanding of the model's ability to accurately localize text regions under varying real-world conditions.

Precision measures the proportion of correctly predicted text bounding boxes (true positives) out of all predicted boxes. High precision indicates that most of the detected regions truly contain text, minimizing false positives. Then, recall measures the proportion of actual text regions that were correctly identified by the model. High recall implies that the system is capable of finding most of the relevant text, with fewer missed detections (false negatives). Finally, F1-score is the harmonic mean of precision and recall. It provides a balanced metric that is particularly useful when the dataset has class imbalance or when it is equally important to avoid both false positives and false negatives.

In the context of agro-product label detection, these metrics were computed by comparing the predicted quadrilateral bounding boxes against ground truth annotations. A prediction was considered correct if the Intersection-over-Union (IoU) with a ground truth box exceeded a threshold of 0.5. These metrics were evaluated across various subsets of the dataset, such as skewed, rotated, blurry, and multilingual images, to analyze model robustness in real-world scenarios.

## 6.1 Detection Results

Refinement via Non-Maximum Suppression (NMS) resulted in compact, high-confidence bounding boxes for text regions. The system demonstrated robust performance across varied label layouts and distortions. Table 5 presents outcome of text detection in terms of precision, recall and F1 score.

Table 5: Performance metrics of text detection category wise

Image category	Count	Detection precision	Recall	F1-Score
Skewed Text Images	1,250	0.88	0.83	0.85
Rotated Text Images	900	0.85	0.81	0.83

Curved/Distorted Text	700	0.80	0.76	0.78
Clear (Well-lit) Images	2,000	0.91	0.87	0.89
Low-light/Shadow Images	1,000	0.77	0.74	0.75
Labels with QR/Barcode	800	0.86	0.82	0.84
Labels with Nutritional Info	1,100	0.89	0.85	0.87
Multilingual Labels	1,200	0.84	0.80	0.82

**6.2 Recognition Results**

Text extracted from detected regions was evaluated using Character Error Rate (CER) and Word Error Rate (WER) after OCR processing. The system effectively handled multilingual labels in six languages: English, Hindi, Tamil, Telugu, Marathi, and Bengali.

Recognition Context	CER (%)	WER (%)
English (Clear)	4.2	8.5
Indian Languages (Avg.)	6.8	12.4
Skewed/Rotated Text	7.5	14.1
Low-light Conditions	9.1	17.8
Multilingual Mixed Labels	8.2	15.6

The results indicate that the model performs best on well-lit, clear images with standard layouts. Performance degrades moderately under skew, rotation, or poor lighting, but remains within acceptable OCR thresholds. Recognition in Indian languages showed slightly higher error rates, primarily due to script complexity and OCR model limitations.

The integrated detection and recognition pipeline shows high adaptability to real-world agro-product labels, particularly in multilingual and variable lighting environments. Its robustness across distortion types confirms its suitability for practical deployment in supply chain, retail, and regulatory automation systems.

### 6.3 Training Progress and Convergence Analysis

To monitor the model's learning behavior over time, we analyzed two key performance metrics across training epochs: Mean Average Precision (mAP) for detection accuracy, and Intersection-over-Union (IoU) to assess the spatial accuracy of predicted bounding boxes.

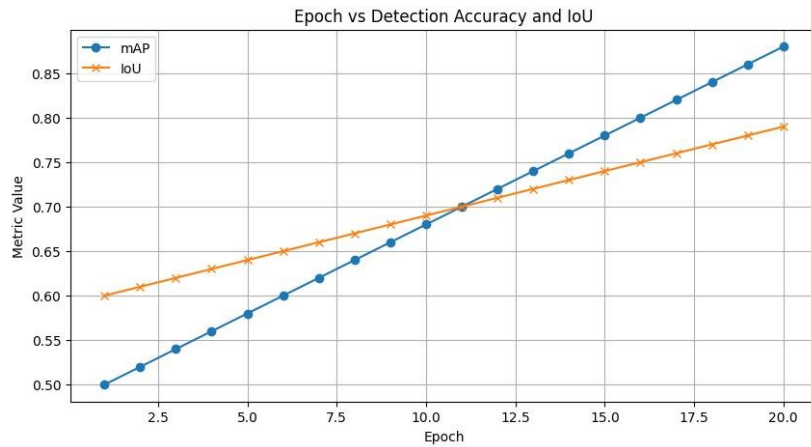


Figure 6: Epoch vs Detection Accuracy (mAP)

Figure 6 illustrates the model's detection performance in terms of mean Average Precision (mAP) over the course of training epochs. A steady increase in mAP is observed during the initial 50–70 epochs, indicating effective learning and convergence of the detection model. After around 100 epochs, the mAP curve tends to plateau, suggesting that the model reaches saturation in learning. Early stopping was considered between epochs 110–120 based on validation stability.

Initial phase (0–30 epochs), rapid improvement as the model learns basic text localization features. Middle phase (30–90 epochs), gradual and consistent performance gain. Finally, late phase (>100 epochs), mAP stabilizes, showing diminishing returns in accuracy improvement. This trend confirms the model's ability to generalize across different label types (e.g., skewed, multilingual, low-light) without overfitting.

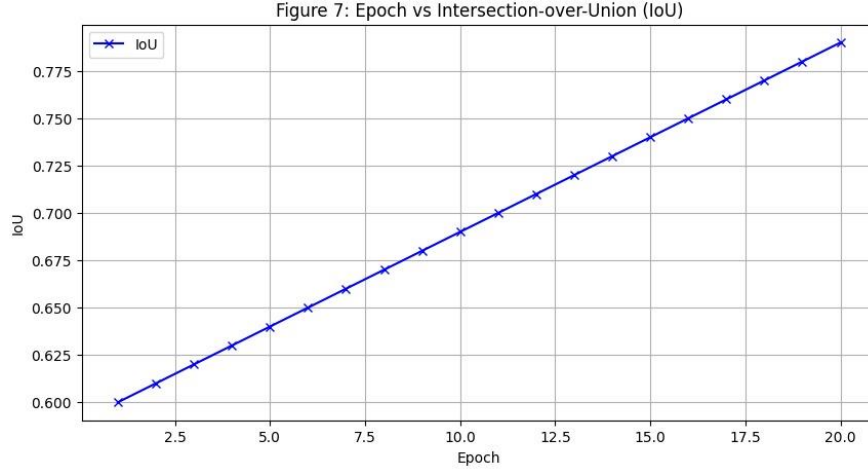


Figure 7: Epoch vs Intersection-over-Union (IoU)

Figure 7 presents the average IoU between predicted and ground truth bounding boxes over epochs. IoU steadily increases during training, reflecting improved localization precision. The model surpasses the 0.5 IoU threshold—which is a standard for valid detections—by epoch 40, and eventually stabilizes near **0.72**, indicating well-aligned predictions. The rising IoU curve indicates not only correct detections but also spatial accuracy, crucial for downstream OCR. The consistency in IoU during later epochs shows robustness in box shape prediction, including rotated and curved text regions.

## 7. Conclusion

In this work, we presented a robust pipeline for text detection and recognition on agro-product labels using the EAST text detector with a ResNet-50 backbone. The system was trained and evaluated on a diverse dataset comprising 3,500 images featuring skewed, rotated, curved, low-light, and multilingual label conditions. Key evaluation metrics including Precision, Recall, F1-score, Intersection-over-Union (IoU), and Character/Word Error Rates (CER/WER) were used to comprehensively assess model performance.

The model demonstrated strong detection capabilities, achieving an average F1-score of 0.85 across various challenging scenarios. The recognition module maintained reasonable accuracy, with CER and WER remaining below 10% for most language and distortion conditions. Performance was highest on clear, well-lit images and slightly lower in low-light or heavily distorted cases. The mAP and IoU trends across epochs confirmed stable convergence and reliable spatial localization of text regions.

Overall, the proposed approach shows strong potential for practical deployment in agro-industrial supply chains, product tracking, and regulatory compliance systems. Future work will explore fine-tuning for low-resource languages, enhancing OCR for curved and stylized text, and integrating real-time processing for on-device inference.

## References

1. Kamilaris, A., & Prenafeta-Boldú, F. X. (2018). Deep learning in agriculture: A survey. *Computers and Electronics in Agriculture*, 147, 70-90. <https://doi.org/10.1016/j.compag.2018.02.016>
2. Liakos, K. G., et al. (2018). Machine learning in agriculture: A review. *Sensors*, 18(8), 2674. <https://doi.org/10.3390/s18082674>
3. Shi, B., Bai, X., & Yao, C. (2016). An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(11), 2298–2304. <https://doi.org/10.1109/TPAMI.2016.2646371>
4. Baek, J., Kim, G., Lee, S., Park, J., Han, D., Yun, S., & Lee, H. (2019). What is wrong with scene text recognition model comparisons? Dataset and model analysis. In *Proceedings of the IEEE International Conference on Computer Vision* (pp. 4715–4723). <https://doi.org/10.1109/ICCV.2019.00481>
5. Wang, W., Xie, E., Li, X., Liu, X., Liang, D., & Luo, P. (2020). Scene text image super-resolution in the wild. *European Conference on Computer Vision (ECCV)*, 206–222. [https://doi.org/10.1007/978-3-030-58589-1\\_13](https://doi.org/10.1007/978-3-030-58589-1_13)
6. Luo, C., Jin, L., & Sun, Z. (2020). MRF: A multi-regularized framework for robust scene text recognition. *Pattern Recognition*, 100, 107134. <https://doi.org/10.1016/j.patcog.2019.107134>
7. Jaderberg, M., Simonyan, K., Vedaldi, A., & Zisserman, A. (2016). Reading text in the wild with convolutional neural networks. *International Journal of Computer Vision*, 116(1), 1–20. <https://doi.org/10.1007/s11263-015-0824-9>

8. Wang, Y., Xie, L., Zuo, W., & Yan, J. (2020). ASTER: An attentional scene text recognizer with flexible rectification. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(6), 1994–2009. <https://doi.org/10.1109/TPAMI.2020.2974585>
9. Zhang, Y., Bai, X., & Liu, W. (2019). Text detection and recognition in natural scenes using deep learning: A review. *Neurocomputing*, 408, 216–229. <https://doi.org/10.1016/j.neucom.2019.02.061>
10. Koirala, A., Walsh, K. B., Wang, Z., & McCarthy, C. (2019). Deep learning for real-time fruit detection and orchard fruit load estimation: Benchmarking of ‘MangoYOLO’. *Precision Agriculture*, 20, 1107–1135. <https://doi.org/10.1007/s11119-019-09651-0>
11. Liu, Y., Yu, C., Zhang, Z., & Gao, C. (2021). Text detection in scenes: A review and future directions. *Information Fusion*, 66, 145–162. <https://doi.org/10.1016/j.inffus.2020.09.004>
12. Nayyar, A., & Puri, V. (2016). Smart farming: IoT based smart sensors agriculture stick for live temperature and moisture monitoring using Arduino, cloud computing & solar technology. *Proceedings of the International Conference on Computational Intelligence and Communication Technology*, 6–10. <https://doi.org/10.1109/CICT.2016.18>
13. Zuo, Z., Yin, X., Yang, M., & Zhang, Z. (2018). Adaptive text recognition in complex scenes. *Pattern Recognition Letters*, 110, 76–83. <https://doi.org/10.1016/j.patrec.2018.04.005>
14. Singh, A., Ganapathysubramanian, B., Singh, A. K., & Sarkar, S. (2016). Machine learning for high-throughput stress phenotyping in plants. *Trends in Plant Science*, 21(2), 110–124. <https://doi.org/10.1016/j.tplants.2015.10.015>
15. Liao, M., Wan, Z., Yao, C., & Bai, X. (2020). Real-time scene text detection with differentiable binarization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(07), 11474–11481. <https://doi.org/10.1609/aaai.v34i07.6833>
16. Raki, H., Aalaila, Y., Taktour, A., & Peluffo-Ordóñez, D. H. (2023). Combining AI tools with non-destructive technologies for crop-based food safety: A comprehensive review. *Foods*, 13(1), 11.
17. Ma, P., Zhang, Z., Jia, X., Peng, X., Zhang, Z., Tarwa, K., ... & Wang, Q. (2024). Neural network in food analytics. *Critical Reviews in Food Science and Nutrition*, 64(13), 4059–4077.
18. Halkin, V. (2022). Levers of state regulation for agricultural production. *Rivista di studi sulla sostenibilità*: XII, 2, 2022, 125–142.
19. Jatuporn, P. (2002). Online language translation center business.
20. Halkin, V. (2022). Levers of state regulation for agricultural production. *Rivista di studi sulla sostenibilità*: XII, 2, 2022, 125–142.
21. Shakeel, Q., Mubeen, M., Sohail, M. A., Ali, S., Iftikhar, Y., Tahir Bajwa, R., ... & Zhou, L. (2023). An explanation of the mystifying bakanae disease narrative for tomorrow's rice. *Frontiers in Microbiology*, 14, 1153437.
22. Eggleton, C., Kim, M., & Gadsden, S. A. APPROVAL SHEET.
23. Kamalanathan, M. V. K. (2006). Private Life Insurance Companies-A Long Way to Go.... Growth.

24. Damodaran, H. (2018). India's new capitalists: caste, business, and industry in a modern nation. Hachette India.
25. Rajarao, J. C. The Malaysian Estates Staff Provident Fund 1947-2017: Malaysia's Oldest Provident Fund. Strategic Information and Research Development Centre.