
Application of the K-Nearest Neighbor (K-NN) Algorithm for Detecting Banana Harvest Feasibility

Citra^{1*}, Arnah Ritonga², Arnita³, Said Iskandar Al Idrus⁴, Debi Yandra Niska⁵

^{1,3,4,5}Medan State University, Faculty of Mathematics and Natural Sciences, Computer Science, Jalan W. Iskandar Psr V Medan Esatate Kab. Deli Serdang, Indonesia

²Medan State University, Faculty of Mathematics and Natural Sciences, Mathematics, Jalan W. Iskandar Psr V Medan Esatate Kab. Deli Serdang, Indonesia

Keywords

GLCM; Harvest Suitability; HSV; K-Nearest Neighbor; Web-based System

*Corresponding Author:

citracitra@mhs.unimed.ac.id

Abstract

This study focuses on detecting banana harvest feasibility at the *green-ripe* stage, an area often overlooked in previous studies that focus only on general ripeness. The objective of this research was to develop a system based on the K-Nearest Neighbor (K-NN) algorithm to classify bananas as “Ready for Harvest” or “Not Ready for Harvest” using digital image processing. The system utilizes Hue Saturation Value (HSV) for color analysis and Gray Level Co-occurrence Matrix (GLCM) for texture identification. Unlike other methods, the combination of HSV and GLCM provides richer, complementary features, improving classification accuracy. The study was conducted at a banana plantation in Kwala Bekala Village, Medan Johor District, with 200 banana images taken from five different locations. The K-NN algorithm, with a value of $K = 3$, was chosen to avoid tie votes and ensure computational efficiency. The system achieved an accuracy of 94%, with precision of 93.5%, recall of 92.8%, and an F1-score of 93%. In beta testing with 33 respondents (18 farmers and 15 non-farmers), the system achieved a user satisfaction rate of 90%. Misclassifications occurred due to factors such as lighting variations and background noise. The study demonstrates the practical benefit of using the K-NN algorithm for determining the optimal harvest time, helping farmers make more accurate decisions, reducing waste, and increasing market competitiveness. This research fills the gap in detecting *green-ripe* bananas, providing an innovative solution to optimize harvest timing in the agricultural industry.

1. Introduction

The development of digital technology in the era of Industry 4.0 has driven transformation in various sectors, including agriculture, through the application of artificial intelligence, big data, and image processing. These technologies enable automatic analysis of crop and environmental conditions, helping farmers make more accurate and timely decisions. In banana production, determining the optimal harvest time is crucial for maintaining fruit quality and competitiveness in the market, especially for export-quality bananas [1].

Indonesia is one of the world's largest banana producers, with production reaching 9.34 million tons in 2023 [2]. However, improper harvest timing continues to be a challenge, leading to high levels of spoilage. Harvesting too early results in unripe fruit, while delayed harvesting causes degradation in physical quality, which can significantly lower market value, particularly in export markets with strict quality standards [3], [4].

Based on an observational study at the Kwala Bekala banana plantation in Medan, the determination of harvest suitability is still done manually and subjectively, leading to inconsistencies and potential classification errors [3]. To address these issues, digital image processing can be employed to detect visual characteristics of bananas based on color and texture, providing a more objective and accurate assessment.

The Hue Saturation Value (HSV) method has proven effective in representing color as perceived by the human eye, while the Gray Level Co-occurrence Matrix (GLCM) method analyzes the texture characteristics of the fruit's surface. The combination of HSV and GLCM is believed to enhance classification accuracy compared to using a single method alone. Unlike methods such as Local Binary Patterns (LBP) or Principal Component Analysis (PCA), which primarily focus on either texture or dimensionality reduction, HSV + GLCM offers a richer combination of color and texture that is more closely aligned with human visual perception. This combination has been shown in previous studies to improve the overall classification performance, particularly when distinguishing between subtle variations in fruit ripeness and quality [5].

The K-Nearest Neighbor (K-NN) algorithm was selected for this study due to its simplicity, effectiveness on small datasets, and lack of need for explicit training. Previous studies have shown that K-NN excels in classifying fruit ripeness, such as bananas and mangoes [6], [7]. However, its specific application for detecting banana harvest readiness, particularly at the green-ripe phase, remains limited.

This study aims to apply a classification system that detects the harvest readiness of bananas by combining HSV color features, GLCM texture features, and the K-NN algorithm. The system is expected to help farmers determine the harvest time more accurately, reducing waste, improving market competitiveness, and ensuring better quality in banana production.

Bananas (*Musa spp.*) are an important horticultural commodity in Indonesia and worldwide. As a climacteric fruit, the ripening process continues after harvesting. Therefore, ideal harvesting is done at the green ripe stage, when the fruit's angle disappears and the skin color gradually turns pale green [8]. The selection of the appropriate harvest time affects fruit quality, shelf life, and export suitability. Under-ripe bananas tend not to develop properly, while overripe bananas have a risk of damage and reduced market value [9].

An image is the result of capturing light on a two-dimensional surface. Digital images or pictures are currently very popular. From a mathematical perspective, an image is a continuous function of light intensity on a two-dimensional surface. A light source illuminates an object, which reflects back part of the light beam. This reflected light is captured by optical devices, such as the human eye, cameras, scanners, and so on, thereby recording the object's shadow, which is referred to as an image. Images as the output of a data recording system can be categorized into several types, namely: optical images in the form of photographs, analog images in the form of video signals as displayed on a television monitor, and digital images that can be directly stored on storage media, such as magnetic tape. The term digital image processing generally refers to the process of manipulating two-dimensional images using a computer. This process involves various steps, such as adjusting gradation, modifying, and manipulating images to achieve specific objectives [10].

Feature extraction is the process of finding unique characteristic values possessed by an image. With these values, an image can be distinguished from other images. The characteristic values produced through the feature extraction process are then used as support vectors [11]. The main purpose of feature extraction is to reduce the data dimension and present the image in a simpler and more meaningful representation [12].

HSV is a color model consisting of three components: Hue (range 0° – 360°), Saturation (saturation level), and Value (brightness) [13]. HSV is more representative of the human visual system and is effective for visual

classification of objects such as fruits [14]. The HSV histogram is used to measure color distribution in an image to support the classification process.

GLCM is a statistical technique for analyzing image texture by mapping the frequency of occurrence of pixel value pairs in specific directions (0°, 45°, 90°, and 135°). GLCM features are used to recognize object surface characteristics, such as roughness, uniformity, and density [15]. The combination of color and texture features has proven to be more effective than single use in various classification studies [5].

K-Nearest Neighbor (K-NN) is categorized as a non-parametric classification algorithm that determines the class of data based on its proximity to a number of nearest neighbors (k). No explicit training process is required, making this algorithm efficient for small data sets [16]. The distance between data points is calculated using Euclidean distance. K-NN is known for its simplicity, flexibility, and high accuracy in visual feature-based classification [17], [18]. The mathematical equation for calculating Euclidean distance in the K-NN algorithm is presented as follows:

$$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2} \quad (1)$$

Explanation:

$d(x, y)$ = Distance between two points x and y

x_i = Feature value i of point x

y_i = Feature value i of point y

n = Number of features

Related research has been conducted on banana ripeness classification using digital image processing methods. In [6], the HSV method and K-NN algorithm were used to classify banana ripeness, achieving an accuracy of 93.33%. Another study in [19] combined GLCM texture features and the K-NN algorithm to distinguish between ripe and unripe bananas, achieving an accuracy of 90%. However, most of these studies are still limited to general ripeness classification and have not specifically targeted harvest suitability detection in the green-ripe phase, so a new, more specific approach is needed.

2. Research Method

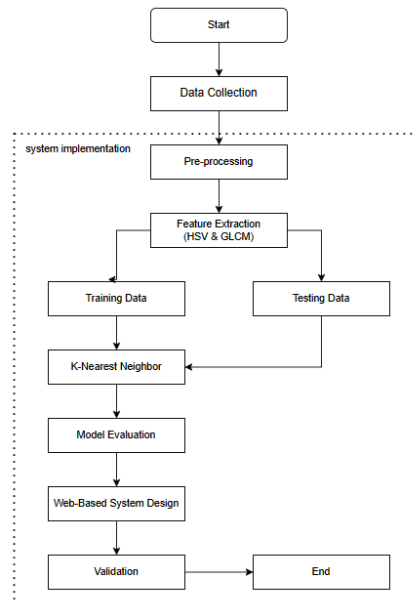


Figure 1. Research Flow Chart

The flowchart in Figure 1 illustrates the main stages in the process of classifying banana harvest suitability using a combination of HSV, GLCM, and K-Nearest Neighbor (K-NN) algorithms. The process is divided into three main parts: data collection, system implementation, and final validation. The process begins with the collection of data in the form of 200 images of bananas in bunches, divided evenly into 100 images of bananas in bunches that are ready for harvest and 100 images of bananas in bunches that are not yet ready for harvest.

The pre-processing stage is crucial to prepare the data for training. It involves several steps, such as background removal, cropping, conversion to grayscale, RGB to HSV transformation, and data augmentation. Background removal ensures that only the banana fruit is focused on in the image, while cropping isolates the fruit to eliminate any unnecessary information. Grayscale conversion simplifies the image by reducing it to a single intensity channel, which is particularly useful for extracting texture features. The RGB to HSV transformation is applied because HSV color space better represents color as perceived by humans, making it more suitable for distinguishing ripeness levels. Data augmentation techniques such as rotation, flipping, zooming, and brightness adjustment are used to simulate real-world conditions like lighting and angle variations that can occur during field photography. This not only increases the dataset's size but also enhances the model's ability to generalize and reduce the risk of overfitting. The final dataset after augmentation consists of 800 images, equally divided between the two classes, *Ready for Harvest* and *Not Ready for Harvest*.

After the images are processed, two types of features are extracted. Color features from the HSV (Hue, Saturation, Value) space. Texture features using the Gray Level Co-occurrence Matrix (GLCM) method, with parameters such as contrast, correlation, energy, and homogeneity. Data division was performed in three ratio scenarios (90:10, 80:20, and 70:30) to evaluate the effect of the training and testing data proportions on model performance. Data randomization in each scenario is implemented using the random state parameter.

Test data is classified based on proximity to training data using the K-Nearest Neighbor (K-NN) algorithm. Euclidean distance is used to calculate the proximity between data points [20]. The Euclidean distance is used because it is simple, efficient, and works well in small datasets like the one in this study. The K-NN algorithm is sensitive to the value of K, which determines the number of nearest neighbors considered for classification. In this study, values of K= 3,5,7,9 are chosen because odd values help avoid ties in classification decisions, ensuring that the system always selects a single class. Furthermore, these values provide a range of small to large neighbor sets, offering a balance between high accuracy and computational efficiency.

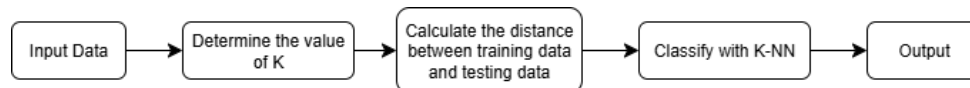


Figure 2. K-NN FlowChart

The classification results are analyzed using an evaluation matrix, namely accuracy, precision, recall, and F1-score, to measure the system's performance on the test data. After the classification model was proven effective, the system was implemented in the form of a web-based application using the Flask framework. This system allows users (farmers) to upload images and obtain classification results automatically. The final stage is system validation through beta testing, where users test the application and provide feedback through a questionnaire. The process concludes with conclusions from the validation results.

3. Result and Discussions

The feature extraction process was carried out using two main approaches, namely color extraction using the Hue Saturation Value (HSV) model and texture extraction using the Gray Level Co-occurrence Matrix (GLCM). The output of this process was numerical values representing each feature, which were used as input for the K-Nearest Neighbor (K-NN) classification model.

Color feature extraction is performed by converting the RGB image into the HSV color space to obtain the Hue, Saturation, and Value values. Next, the mean values of these three components are calculated in the area that

has been masked as a banana. Figure 3 shows the difference in color component distribution between the “Ready for Harvest” and “Not Ready for Harvest” classes.

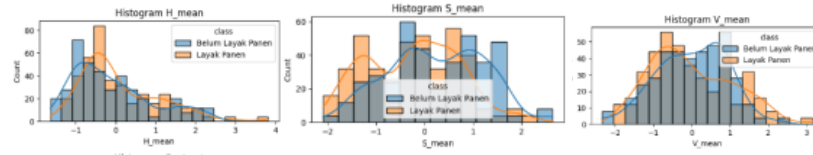


Figure 3. Visualization of HSV Distribution

In Figure 3, it can be seen that the distribution of Hue values shows a fairly clear separation between the two classes, where the Harvestable class tends to have higher Hue values. In the Saturation and Value components, there is overlap, although there is a tendency for the Not Harvestable class to have higher Value values.

Texture features were extracted from the grayscale image obtained through masking using the GLCM matrix with an angle parameter of 0° and a distance of 1 pixel. The resulting features include Contrast, Correlation, Energy, and Homogeneity. The visualization of the GLCM feature distribution is presented in Figure 4.

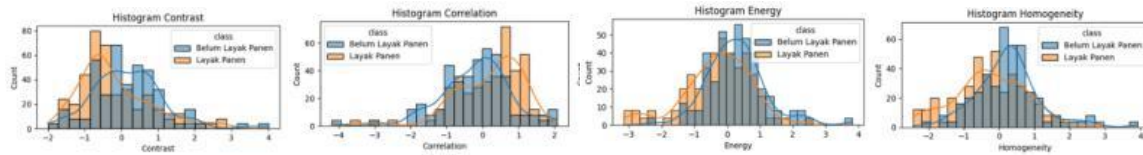


Figure 4. Visualization of GLCM Distribution

Image classification is performed using the K-Nearest Neighbor (K-NN) algorithm, utilizing the extracted HSV and GLCM feature values as input. The distance between the test data and the training data is calculated using the Euclidean formula to determine the similarity between samples.

Experiments on the value of K were also conducted in addition to data partitioning. In this study, experiments were conducted with K values of 3, 5, 7, and 9. The results of these experiments are presented in Table 1 below.

Table 1. Comparison of Accuracy in Data Models

Description	70%:30%	80%:20%	90%:10%
K = 3	90%	94%	100%
K = 5	81%	85%	88%
K = 7	78%	76%	78%
K = 9	77%	81%	81%

Based on Table 1, several K parameter values were experimentally tested, namely K = 3, 5, 7, and 9. The model was then tested on three variations of training and test data splits, namely 90%:10%, 80%:20%, and 70%:30%. The test results show that the value K = 3 provides the best classification results across all data division scenarios. The highest accuracy of 100% was achieved in the 90%:10% data division, while the 80%:20% division resulted in an accuracy of 94%.

Conversely, models with K = 7 and K = 9 produced lower accuracy across all tested data splits. For the model with K = 7, the accuracy ranged from 76% to 78%, while the model with K = 9 showed similar accuracy, with values ranging from 77% to 81%. Therefore, the model selected for further analysis is the model with a K value of 3 and an 80%:20% data split, as it provides a balance between high accuracy and the model's ability to generalize well.

The confusion matrix results for both classes can be seen in Figure 4. Each class contains confusion matrix values that can be analyzed for each class category. The following is the confusion matrix for all categories.

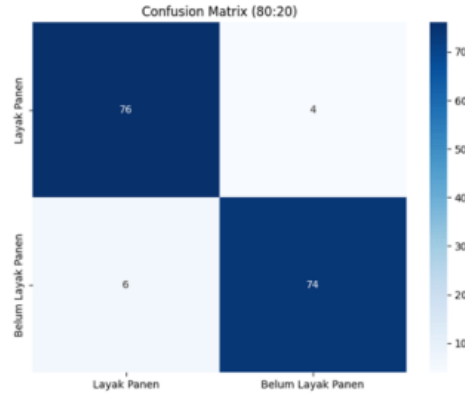


Figure 5. Confusion Matrix

Based on the confusion matrix in Figure 5 above, it can be seen that the model incorrectly predicted 10 data points out of 160 test data points. Four data points with the original label “Ready for Harvest” were predicted as “Not Ready for Harvest,” while six data points with the original label “Not Ready for Harvest” were predicted as “Ready for Harvest.” The accuracy results for the “Ready for Harvest” and “Not Ready for Harvest” categories with a 80%:20% ratio and K = 3 can be calculated using the following formula.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \times 100\%$$

$$Accuracy = \frac{76 + 74}{76 + 74 + 4 + 6} \times 100\% = 94\%$$

The precision, recall, and F1-score values for each class show consistent performance, with average scores above 90%. Cross-validation curve testing was performed to evaluate model stability and detect potential overfitting, as shown in Figure 6 below.

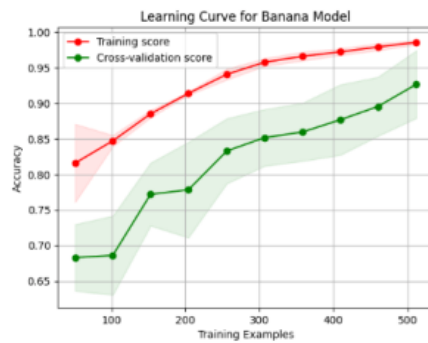


Figure 6. Cross Validation Curve

The 5-fold cross-validation results show the stability of the model with low deviation between the training and test data, so the model is declared reliable in detecting the harvestability of bananas based on images. A web-based banana harvest suitability classification system was developed by combining the K-NN algorithm with HSV and GLCM features, providing instant classification results from uploaded images of bananas in bunches.

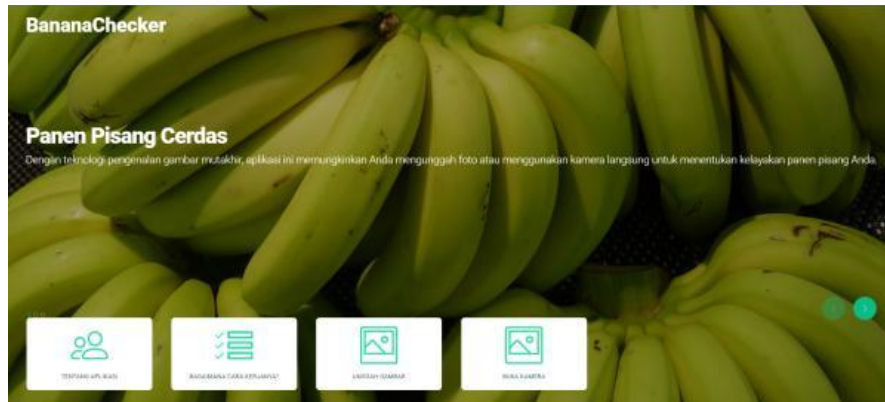


Figure 7. Web-based Banana Checker Application Display on Desktop



Figure 8. Web-based Banana Checker Application Display on Mobile

Figures 7 and 8 show the Banana Checker application interface, designed to be responsive with a vertical layout on mobile and horizontal layout on desktop to ensure user-friendly access across various devices.

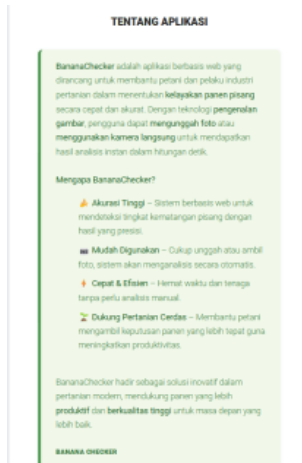


Figure 9. Banana Checker Application Menu Display



Figure 10. Banana Checker Application Functionality Menu Display

Figure 9 shows the Banana Checker application information menu, which includes an explanation of the system's purpose and benefits, designed to help farmers accurately and efficiently determine harvest readiness through image recognition technology. Figure 10 shows the Banana Checker application functionality menu display, which is outlined in four main steps: images are uploaded or captured directly, analyzed quickly and accurately, classified instantly, and the results are used to optimize harvest timing.

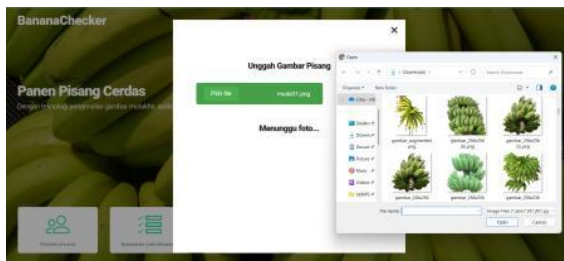


Figure 11. Banana Checker Application Image Upload Menu Display

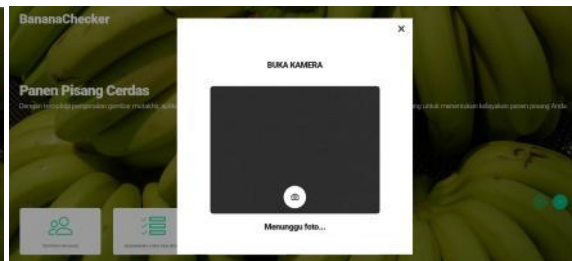


Figure 12. Banana Checker App Camera Open Menu Display

Figure 11 shows the image upload menu display, where banana images can be selected from the device and automatically processed by the system to determine harvest readiness. Figure 12 shows the Camera Open menu designed to enable direct photo capture. Upon access, the system displays the message "Waiting for photo..." as an indicator of readiness to process the image for harvest suitability analysis.



Figure 13. Banana Checker application (after image input)

Figure 13 shows the Banana Checker classification results interface, where the banana image is classified as “Not Ready for Harvest” with a probability of 100%. This value indicates the model's full confidence in the classification result provided.

The system evaluation was conducted through beta testing involving 33 respondents, consisting of 18 farmers and 15 non-farmers. Ten aspects were tested, including ease of use, system stability, responsiveness, classification accuracy, and feature completeness. The assessment was conducted using a Likert scale questionnaire (1–4). The results of the beta testing of the ten questions given to respondents regarding the Banana Checker website show that the total score of all questions was calculated to determine the usability of the Banana Checker website.

Table 2. Quality Criteria

Assessment Percentage	Interpretation
81% - 100%	Very Usable
61% - 80%	Usable
41% - 60%	Sufficiently Suitable
21% - 40%	Less Suitable
0% - 20 %	Not Suitable

The following are the results of the beta test calculations:

$$\begin{aligned}
 \text{Total} &= 95 + 90 + 86 + 89 + 87 + 92 + 89 + 91 + 89 + 90 \\
 &= \frac{898}{10} \times 100\% \\
 &= 90\%
 \end{aligned}$$

Based on these results, it can be concluded that the Banana Checker website is highly suitable for use.

Percentage of Beta Testing Results for the BananaChecker Website

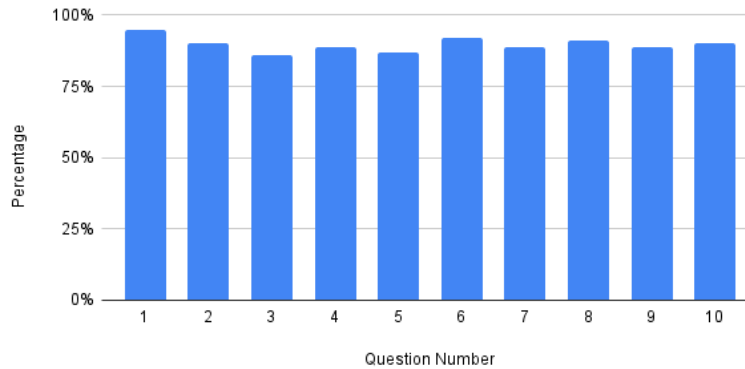


Figure 14. Beta Testing Results Graph

4. Conclusions and Future Works

Based on the research results, it can be concluded that the K-Nearest Neighbor (K-NN) algorithm combined with HSV and GLCM features successfully addresses the research gap identified in the introduction by detecting green-ripe harvest feasibility rather than general ripeness. The best configuration, using an 80:20 data split and K=3, achieved 94% accuracy, demonstrating stable performance with no signs of overfitting and validated through beta testing as highly usable in real farming scenarios. This level of accuracy helps reduce incorrect harvest decisions by less than 10%, improving fruit quality, reducing waste, and increasing market value, particularly for export-quality bananas, thereby offering practical benefits to farmers through the Banana Checker web application. Future work should focus on expanding the dataset and exploring additional feature

extraction methods such as Local Binary Pattern (LBP) or Shape Descriptors to enhance classification robustness and generalization.

5. References

- [1] S. Yasmin and M. Billah, "Digital image processing applications in agriculture with a machine learning approach," *Agric. Sci. Technol.*, vol. 15, no. 4, pp. 12–22, Dec. 2023.
- [2] BPS, "Statistik Hortikultura," p. 32, 2023.
- [3] S. P. Adenugraha, V. Arinal, and D. I. Mulyana, "Klasifikasi Kematangan Buah Pisang Ambon Menggunakan Metode KNN dan PCA Berdasarkan Citra RGB dan HSV," *J. Media Inform. Budidarma*, vol. 6, no. 1, pp. 9–17, Jan. 2022.
- [4] Rifki Kosasih, "Classification of Banana Ripe Level Based on Texture Features and KNN Algorithms," *J. Nas. Tek. Elektro dan Teknol. Inf*, vol. 10, no. 4, pp. 383–388, Nov. 2021.
- [5] W. Shinta Sari and C. Atika Sari, "Klasifikasi Bunga Mawar Menggunakan KNN dan Ekstraksi Fitur GLCM dan HSV," *SKANIKA Sist. Komput. dan Tek. Inform.*, vol. 5, no. 2, pp. 145–156, 2022.
- [6] Lia Kamelia, Mufid Ridlo Effendi, and Nisa Hafidzatul Adila, "Ripeness Level Classification of Banana Fruit Based on Hue Saturate Value (HSV) Color Space Using K-Nearest Neighbor Algorithm," *Int. J. Adv. Trends Comput. Sci. Eng.*, vol. 10, no. 2, pp. 941–947, Apr. 2021.
- [7] Mutmainnah Muchtar and R. A. Muchtar, "Perbandingan Metode Knn Dan Svm Dalam Klasifikasi Kematangan Buah Mangga Berdasarkan Citra Hsv Dan Fitur Statistik," *J. Inform. dan Tek. Elektro Terap.*, vol. 12, no. 2, pp. 876–884, Apr. 2024.
- [8] D. Surya Prabha and J. Satheesh Kumar, "Assessment of banana fruit maturity by image processing technique," *J. Food Sci. Technol.*, vol. 52, no. 3, pp. 1316–1327, 2015.
- [9] S. Maduwanthi and R. Marapana, "Biochemical changes during ripening of banana: A review," *Int. J. Food Sci. Nutr.*, vol. 2, no. 5, pp. 166–169, Sep. 2017.
- [10] S. Irianto, "Pengolahan Citra Digital," Bandar Lampung, 2014.
- [11] Sumijan and Pradani Ayu Widya Purnama, *Teori dan Aplikasi Pengolahan Citra Digital Penerapan dalam Bidang Citra Medis*, Cetakan pe. CV Insan Cendekia Mandiri, 2021.
- [12] F. P. Rhesal Mahadyanto, Diah Arifah Prastiningtyas, "Penerapan Metode Jaringan Syaraf Tiruan Radial Basis Function untuk Identifikasi Jenis Mangga Berdasarkan Pola Daun," *J-INTECH*, vol. 07, no. 1, pp. 90–96, 2019.
- [13] G. P. A. Saputra, "Analysis of Coffee Bean Roasting Maturity Levels Using Color Feature Extraction," *J-INTECH*, vol. 12, no. 1, pp. 123–128, 2024.
- [14] V. Chernov, J. Alander, and V. Bochko, "Integer-based accurate conversion between RGB and HSV color spaces," *Comput. Electr. Eng.*, vol. 46, pp. 328–337, 2015.
- [15] P. N. A. T. S. Muljono, *Pengolahan Citra Digital*. Yogyakarta : ANDI, 2017.
- [16] Rahmadden, Denok Wulandari, Masni Renova, Gilang Ramadhan, and Retno Sari, *Machine Learning*. 2024.
- [17] Danny Erry Trihandika, "Sistem Pengelolaan Informasi Pertanian Menggunakan Metode Case Based Reasoning pada Gapoktan Sidomakmur," *J-INTECH*, vol. 4, no. 1, pp. 66–70, 2017.
- [18] Arnita, Faridawaty Marpaung, Fitrahuda Aulia, Nita Suryani, and Rinjani Cyra Nabila, *Computer Vision Dan Pengolahan Citra Digital*, Edisi Pert. Pustaka Aksara, 2022.
- [19] E. S. K. D. H. Y. N. I. Sembiring, "Enhanced Banana Ripeness Detection using GLCM and K-NN Methods," 2024 3rd Int. Conf. Creat. Commun. Innov. Technol., pp. 1–7, 2024.
- [20] M. K. Yosua Kristanto, Diah Arifah Prastiningtyas, "Penerapan Algoritma Nearest Neighbor untuk Menentukan Rekomendasi Solusi Terhadap Layanan Kantor Teknologi Informasi STIKI Malang," *J-INTECH*, vol. 7, no. 1, pp. 72–79, 2019.