# Rapid Predicting of Ripeness and Lycopene Content of Fresh Tomatoes using Smartphone Camera Images and Machine Learning

1st Phuong Mai Dinh
*2HUS High School For Gifted Students*
*VNU University of Science, Vietnam*
*National University, Hanoi*
Ha Noi, Viet Nam
dmphuong021109@gmail.com

2nd Khanh Nguyen Ngoc Le
*2HUS High School For Gifted Students*
*VNU University of Science, Vietnam*
*National University, Hanoi*
Ha Noi, Viet Nam
khanhthanh3679@gmail.com

3rd Lam Duc Dinh
*3School of Information and*
*Communication Technology*
*Hanoi University of Science and*
*Technology*
Ha Noi, Vietnam
dinhduclam2001@gmail.com

4th Mai Thi Tuyet Phan*
*1VNU University of Science*
*Vietnam National University, Hanoi*
Hanoi, Vietnam
phanthituyetmai@hus.edu.vn

5th Vu Giang Tran
*4HaNoi-Amsterdam High School for*
*Gifted, Viet Nam*
Ha Noi, Viet Nam
trangiangvu13042009@gmail.com

*Abstract*—**This study developed a fresh tomato ripeness and lycopene content prediction system based on RGB color indices extracted from digital images. An image dataset was created using machine learning to localize and label into 6 different ripeness levels from 1,200 tomato images acquired by a smartphone camera. Next, the subject of the tomato image was extracted from the background using an image GrabCut algorithm. Then, the RGB color index was extracted from the background-separated tomato images. At the same time, the lycopene content of 100 tomatoes in the image data set was determined by UV-VIS spectrophotometry. Then, the correlation models between the RGB color index extracted from the image and the lycopene content were established. Based on established models, a web-based ripeness detection and recognition system to estimate the lycopene content of fresh tomatoes was developed. The results showed that the lycopene content estimated from this system with input images captured by smartphone camera using artificial intelligence (AI) achieved high accuracy and repeatability, allowing the development of a rapid, non-invasive/non-destructive, and in-situ tool to proactively assess the quality of fresh tomatoes before harvest.**

*Keywords*— **Ripeness, Lycopene content, Prediction system, RGB color index, Machine learning.**

## I. INTRODUCTION

Tomato is a popular and highly demanded vegetable crop, second only to potato in global crop production, with an annual production of over 1.8 million tons, and was selected as the target for this study. Color and shape are the main characteristics used to grade or evaluate tomato fruit quality. Tomatoes are famous for their vibrant red color, which not only indicates ripeness and desirable flavor but also indicates relative lycopene content; the red color is mainly due to lycopene, which ranges from yellow to red [1]. Therefore, the prediction of ripeness based on the red color of tomatoes has been evaluated by growers based on visual experience as well as using colorimetric devices. Tomatoes with deep red color are usually riper than those with light red, pink, or orange and green colors, with sweet flavor and high lycopene content [2,3]. Another method is spectral analysis, measuring the reflectance of the surface at specific wavelengths using a spectrometer [4-6]. Although the spectral analysis gives highly accurate results, its implementation is difficult due to the investment in expensive spectrometers and spectrum analyzers, as well as the need for calibration for initial installation and complicated operation.

In recent years, AI and machine learning-based harvest prediction and non-destructive quality parameter estimation technologies have attracted attention, especially in tomato cultivation [7-10]. To predict ripeness, it is necessary to first detect tomatoes in camera images. Object detection methods in images include edge detection, pattern matching, and machine learning, which mainly rely on sliding window algorithms that follow a three-stage process: 1) Region selection, 2) Feature extraction, and 3) Classification task performance. Machine learning includes Region-based Convolutional Neural Network (RCNN), Single-Box Multi-Detector (SSD) [10], and YOLO. Among them, YOLO is currently the most popular method with a relatively fast detection speed and high detection accuracy [11]. Various methods have been used in technology and research related to tomato harvest prediction using spectral images taken by a spectral camera for machine learning [12-15], and Mask-CNN (Convolutional Neural Network) machine learning to accurately identify RGB differences [13]. This technology analyzes greenhouse tomato image data to predict ripeness, yield, and harvest time, allowing farmers to easily plan their supply, preventing oversupply and excess inventory.

Among the nutritional compounds present in tomatoes, lycopene is recognized as the most bioavailable carotenoid, offering both health-promoting and aesthetic benefits. Consequently, lycopene content is a key indicator for assessing tomato quality. However, accurately determining lycopene content remains challenging, particularly for small-scale farmers. Conventional analytical methods require advanced instrumentation, complex procedures, and high operational costs, making them impractical for widespread use in agricultural production. To address this limitation, the present study proposes a prediction model for estimating lycopene content in fresh tomatoes using digital images captured by a smartphone camera. The objective is to develop a rapid, non-destructive approach for determining both lycopene concentration and ripeness based on RGB color data extracted from tomato images acquired under field conditions.

## II. Materials and methods

### 1. Plant material and sample preparation

The tomato cultivar used in this study was "Beef Da Lat", which is one of the most popular tomato cultivars in Vietnamese markets. The tomatoes were cultivated on the Gelponic Greenhouse Farm at Thanh Oai, Ha Noi, Viet Nam. The primary data set of fresh tomato images was acquired using a Galaxy A13 smartphone camera at different stages of ripening, in a real-life cultivation condition (Fig 1).

At the same time, the tomato fruits were harvested at different degrees of ripeness based on fruit size and fruit color. This method yielded a batch of fruit samples with different levels of lycopene content. The fruits were carefully handled during harvest. Samples were screened for differently colored fruit surfaces in a total of 100 selected fruits.



Fig 1. Tomato images at different stages of ripeness in the greenhouse farm.

### 2. Dataset preparation

#### 2.1. Data acquisition

The dataset was created using a smartphone camera to capture 1,200 images of 1152x2560 pixels. To enhance the model's generalization and robustness, factors that influence prediction performance—such as occlusion, shadows, and lighting variability—were carefully considered during data collection. Particular emphasis was placed on illumination and occlusion. To simulate diverse lighting conditions, images were captured at multiple times of day (morning, noon, afternoon, and evening), enabling the model to better adapt to changes in ambient light. Additionally, because tomato fruits are frequently partially obscured by leaves, branches, or neighboring fruits under real cultivation conditions, deliberate occlusion scenarios were introduced. These measures ensured that the model could effectively handle common sources of interference in practical settings, thereby improving its robustness and predictive accuracy. A specific example of data collection is shown in Fig 1. Image size normalization, light balance, color correction, and noise removal were performed to ensure input data consistency.

#### 2.2. Data annotation

In this study, 1,200 tomato image datasets were collected and screened in a real production environment. The whole dataset was divided into a training set, validation set, and test set according to a ratio of 7:2:1.

Labelling software CVAT (Computer Vision Annotation Tool) was used to manually label the tomatoes in the image using bounding boxes. Each bounding box is drawn around a separate tomato fruit, ensuring that the object is completely covered but limiting background or leaves. Since this study constructed the model from the perspective of tomato ripeness grading, a tomato ripeness grading model was constructed into six types of labels: green, breaker, turning, pink, light red, and red.

#### 2.3. Standard YOLOv8 Model

For tomato fruit detection, we employed the YOLOv8 model due to its high real-time detection accuracy and low computational cost. A comprehensive methodology for tomato fruit detection, segmentation, and extraction is discussed by integrating the YOLOv8 model for object detection with Google Colab (GC) as the computational environment for segmentation tasks. The methodology of the proposed approach is shown in Fig 2. The labeled dataset is used to train the YOLO, which enables the system to automatically detect and localize tomatoes in annotated images. The boxes were automatically encased around the tomato fruits themselves by using YOLOv8, which was later refined through GC. The output of the model is the bounding box coordinates along with the corresponding ripeness prediction label for each tomato.
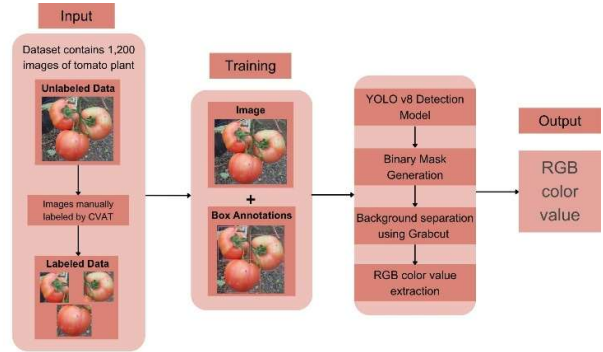


Fig 2. Methodology for tomato fruit detection, segmentation, and RGB extraction.

#### 2.4. Background separation method

Based on the bounding box coordinates generated by the YOLO model, the tomato regions were isolated from the background. The background was completely removed to ensure that only the tomato area remained, enabling more accurate color analysis. From there, a binary mask is created for each fruit, in which:

- Pixels belonging to the tomato region have a value of 1
- Background pixels or irrelevant parts have a value of 0

These masks were then used for subsequent processing steps, including RGB index calculation, lycopene estimation, and more accurate ripeness analysis. Background separation in tomato images was performed using a semi-automatic approach based on a Gaussian Mixture Model (GMM) combined with the GrabCut algorithm to extract RGB color-index features. GrabCut identifies the boundary between the object (foreground) and the background by modeling their color distributions. The resulting image contains a fully isolated tomato, prepared for further processing with the GrabCut algorithm (Fig 2).

#### 2.5. RGB feature extraction

Calculate the average value of three R, G, and B color indices from the pixels in the extracted tomato area (Fig 2).

### 3. Lycopene quantification

#### 3.1. Lycopene quantification

Lycopene content in tomatoes was quantified spectrophotometrically with a UV-Vis spectroscopy apparatus (DV-8200 Single Beam Visible Spectrophotometer-Drawell). For each fruit, the pulp obtained after removing the skin and seeds was used for color measurement. The sample was crushed by machine and then extracted using a 50/50 volume fraction of ethyl acetate and acetone solvent mixture under the

condition of a 2/1 solvent volume/tomato sample mass ratio, mL/g, at 40ºC for 6 h. The extract was then separated from the mixture by a vacuum filter funnel using Advantec filter paper No. 2 (pore size 5 μm, diameter 90 mm). Next, the lycopene in the extract was separated by n-Hexane solvent in a separating funnel (Fig 3). Finally, the lycopene concentration in the n-Hexane solvent of tomato samples was measured using a UV-Vis spectroscopy at 503 nm wavelength [7]. The lycopene content in fresh tomatoes was determined by (1):

$$C = \frac{A_{503}}{a_{503}} \cdot \frac{V}{1000} \cdot \frac{100}{w}, \left(\frac{mg}{100g}\right) \quad (1)$$

Where C is the lycopene content in 100 g of fresh tomatoes (mg/100 g), $A_{503}$ is the absorbance at 503nm wavelength, $a_{503}$ is the specific absorption coefficient of lycopene in n-Hexane, V is the volume of extract (mL), and w is the weight of tomatoes (g).
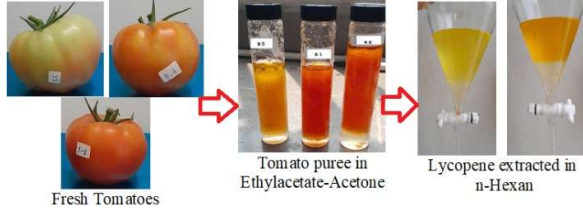


Fig 3. Experimental steps for lycopene extraction from fresh tomatoes.

### 3.2. Lycopene estimation model development

Spectrophotometrically measured lycopene content was compared with fruit RGB values and their derived functions. RGB values were plotted against lycopene to assess relationships, and regression models were fitted. The model with the highest predictive accuracy was incorporated into the application to estimate lycopene content in tomatoes. Using a simple linear regression model with (2):

$$y = a.x + b \quad (2)$$

In which y is the lycopene content (mg/100g fresh fruit), x is the color index value.

Training data is taken from tomato fruit samples that have been analyzed for actual lycopene in the laboratory to ensure the accuracy of the model.

### 4. Application development

With the growing role of digital transformation and AI in smart agriculture, computer vision has become an essential tool for monitoring and assessing crop quality. This project developed a web-based system that enables users to upload tomato images, automatically determine ripeness, estimate lycopene content, and predict optimal harvest time without specialized equipment. The system uses the YOLO model to detect and classify tomatoes by ripening stage. GrabCut and RGB extraction are then applied to isolate the fruit and compute average R, G, and B values. Lycopene content is estimated using the study's empirically derived correlation between color indices and lycopene.

The proportion of correctly detected tomatoes in an image was defined as the accuracy rate, calculated as (3). The detection rate was defined as the ratio of detected tomatoes to the actual number present in the image, expressed as (4).

$$\text{Accuracy rate} = \frac{TP}{TP + FP} \quad (3)$$

where TP and FP denote true and false positives, respectively.

$$\text{Detection rate} = \frac{TP}{TP + FN} \quad (4)$$

where TP and FN denote true positives and false negatives, respectively.

### 4.2. Model evaluation

Model performance for classifying tomato ripeness and defects was evaluated using accuracy, precision, recall, and F1-score. Accuracy measures overall correct classifications but can be misleading with imbalanced data. Precision reflects the proportion of true positives, minimizing false positives, while recall measures the ability to detect actual positives, critical when missing defective fruits is costly. The F1-score balances precision and recall, particularly for imbalanced datasets. However, these metrics do not capture defect severity or subtle ripeness variations, suggesting that complementary or domain-specific measures may be needed for a more complete evaluation.

## III. RESULTS AND DISCUSSIONS

### 1. Create a dataset for machine learning models

Processed images, assigned labels, separated backgrounds and extracted color indices in the RGB color space of 1,200 photos of tomato fruits in greenhouse farm to obtain a database of labeled images into 6 labels corresponding to 6 levels of ripeness (with label names) including: green (green), unripe (Breaker), young (Turning), beginning to ripen (Pink), medium rip (Light red) and overripe (Red) exported in YOLO format (including .txt file containing bounding box coordinates and class labels) used for training and testing. TABLE 1 shows the classification of 6 levels of tomato ripeness with label names, descriptions of color status and ripeness level, and harvest time determined by sensory and experimental evaluation.

TABLE 1. CLASSIFICATION OF FRESH TOMATO LABELS ACCORDING TO 6 LEVELS OF RIPENESS.

| Label | Green | Breaker | Turning | Pink | Light Red | Red |
|---|---|---|---|---|---|---|
| Colour | Green | Light green, light pink | Light pink, light yellow | Dark pink, yellow | Light red, light yellow | Completely Red |
| Ripeness | 0%-20% | 20%-40% | 40%-50% | 50%-70% | 70%-85% | 85-100% |
| Images | | | | | | |

The specific data distributions are shown in TABLE 2. This labeled dataset is further provided to YOLOv8.

TABLE 2. CLASS DISTRIBUTION ACROSS DATASETS.

| Type | Training Set | Validation Set | Test Set | Total | Percentage, % |
|---|---|---|---|---|---|
| Green | 42 | 12 | 6 | 60 | 5 |
| Breaker | 84 | 24 | 12 | 120 | 10 |
| Turning | 84 | 24 | 12 | 120 | 10 |
| Pink | 210 | 60 | 30 | 300 | 25 |
| LightRed | 210 | 60 | 30 | 300 | 25 |
| Red | 210 | 60 | 30 | 300 | 25 |

Note: The "Percentage" column indicates the proportion of each class in the entire dataset.

## 2. Establish a correlation model between RGB color values and lycopene content

This study aimed to non-destructively estimate lycopene content in intact tomato fruits using R, G, and B color indices. Tomatoes at different ripening stages were photographed with a mobile phone to obtain RGB values, and their lycopene content was measured using UV–VIS spectroscopy.

### 2.1. RGB color values and lycopene content

Fig 4 summarizes the measured RGB color values and lycopene content of 48 fresh tomatoes. Due to the different ripening stages, significant differences in R, G, and B color values and lycopene content between samples were observed.

It can be seen that the R (red) color value in the RGB color space ranges from 132 to 197. The G (green) color value in the RGB color space ranges from 49 to 200. The B (blue) color value in the RGB color space ranges from 25 to 148. The lycopene content ranges from 0 to 10.98 mg/100 g fresh tomato. This finding suggests that the R, G, and B color values in the RGB color space can provide useful information on lycopene content in tomato fruit. Notably, the R, G, and B color values all appeared in tomato samples with ripeness varying from green to ripe, but with clearly different value ranges. This shows that the ripeness of the tomato will be reflected by the RGB color index value combined from the individual values, not just each individual color value.
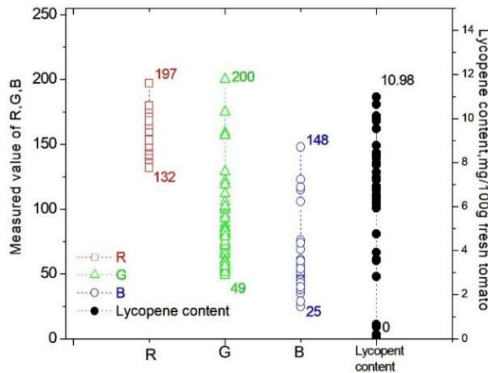


Fig 4. The distribution of the measured RGB color values and lycopene content of fresh tomatoes.

### 2.2. Establish the lycopene content estimation model

Using RGB indices extracted from images of 48 tomatoes at different ripening stages, a correlation model was developed to relate lycopene content to transformed RGB values. Three chromaticity parameters—R/(R+G+B), (2R–G–B)/(R+G+B), and (3R–G–B)/(R+G+B)—were examined. Among these, (2R–G–B)/(R+G+B) provided the best linear fit, yielding an $R^2$ of 0.9769 (Fig. 5). Because lycopene is a red pigment, transformations that enhance red and suppress green and blue intensities are inherently more effective predictors. The (2R–G–B)/(R+G+B) parameter not only showed the highest $R^2$ value but also produced a regression line passing through the origin, further confirming its suitability as a visual predictor of lycopene content in tomato fruits. Ye et al. [12] similarly used color-based modeling by converting RGB values to the CIELab space, where the transformed parameter $(a/b)^2$ achieved the best linear fit with an $R^2$ of 0.81. Overall, the proposed model demonstrated improved accuracy for estimating lycopene content in fresh tomatoes under real cultivation conditions.
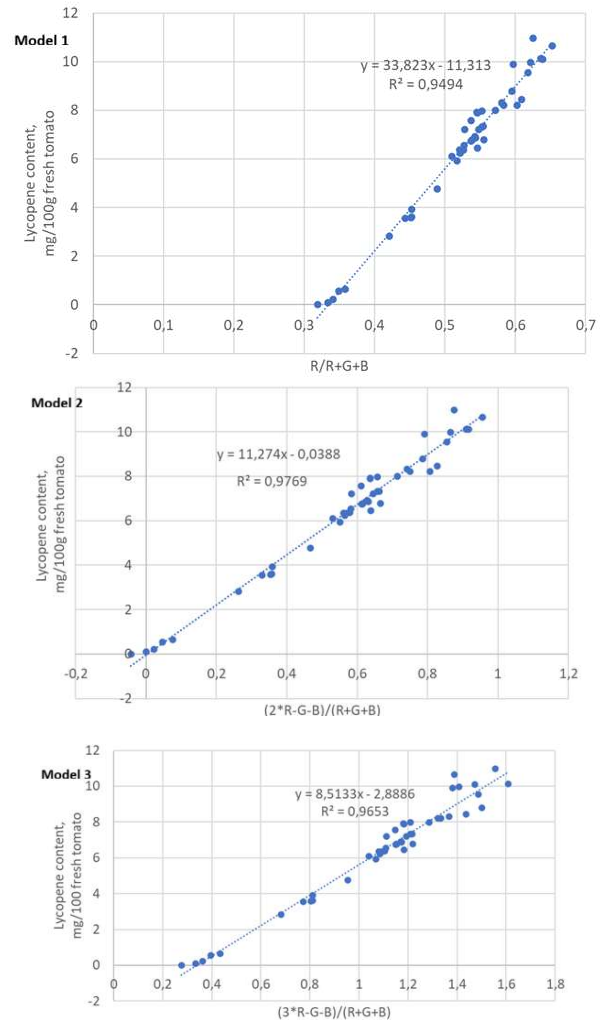


Fig 5. The relationships between lycopene content and the transformed RGB values.

## 3. Development of the application system

A web model has been established to predict the ripeness and lycopene content of tomatoes from images taken by a smartphone camera. Fig 6 shows the graphical user interface for simulating and tracking model performance. The web interface is designed to be user-friendly and intuitive, allowing users to easily download images, view analysis results, color index, lycopene value, and predict harvest time immediately. On this website, users only need to press (1) to download images or press (2) to take a photo directly, and then press (3) to analyze the data. After a quick analysis, the results will be returned, including the number of identified fruits (4), images of the zoned tomatoes (5), and corresponding information boxes (6). In the information box, the model will provide the following contents: the name of the fruit being analyzed (6.1); the RGB index of that fruit (6.2); and the corresponding lycopene content (6.3). The application was evaluated using a test dataset of fresh tomato images captured by a smartphone. The system successfully identified fruits, labeled their ripeness grades, and displayed the selected regions of interest. Estimated lycopene contents were stored locally and could be exported for further analysis, demonstrating the application's effectiveness for lycopene estimation and fruit ripeness.
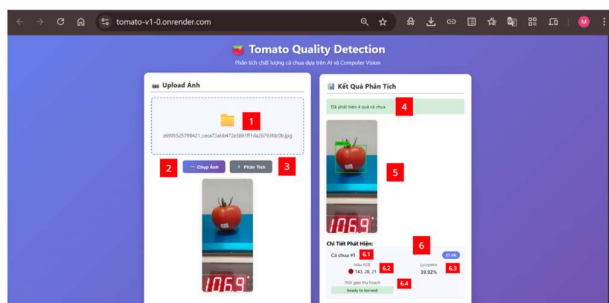
Fig 6. Screen interface of the application system to predict the lycopene content of a fresh tomato image captured using a smartphone camera.

The YOLO-based tomato detection method achieved an accuracy rate of 0.94 and a detection rate of 0.89. This indicates that 94% of detected objects were true tomatoes and that 89% of all tomatoes present in the images were successfully identified, with 11% missed. Because this study did not aim to estimate tomato yield, these omissions do not affect the overall objective. Nonetheless, these results suggest that the model's generalization remains imperfect, which may affect its performance on unseen data.

Ambient lighting can substantially affect camera-acquired RGB values, and variations in mobile camera specifications may also introduce differences in measurements of the same fruit [10, 17, 18]. As a result, slight discrepancies in estimated lycopene content across devices and conditions are unavoidable. Nevertheless, when precise quantification is not required, camera-based imaging remains effective for fruit grading and quality classification. Future work will include evaluating the system's performance across different smartphone cameras and lighting conditions to further assess its robustness in practical applications.

## IV. CONCLUSIONS

We established a comprehensive tomato fruit image dataset consisting of six ripeness classes for machine learning applications, constructed from 1,200 smartphone images collected under real field-farming conditions. Using this dataset, we developed a non-destructive system to estimate lycopene content and ripeness in fresh tomatoes. Our RGB-based regression models identified $(2R–G–B)/(R+G+B)$ as the optimal predictor, achieving an $R^2$ of 0.9769.

The system integrates YOLOv8 for fruit detection and GrabCut for segmentation, enabling rapid lycopene estimation and automated fruit grading from tomato images. The YOLO-based tomato detection method demonstrated high performance, achieving an accuracy of 0.94 and a detection rate of 0.89. Initially implemented as a web application and later as a mobile platform, the system allows users to scan fruits with a smartphone camera and receive instant ripeness assessments.

This developed system is user-friendly and practical, providing real-time detection capabilities. This is particularly valuable in real-cultivations where rapid decision-making is required to determine optimal harvest timing and control tomato quality before production decisions. The system also supports retail use by enabling rapid evaluations and helping consumers assess fruit ripeness before purchase. Our results demonstrate the model's effectiveness, practical impact, and potential for scalable, in-field deployment in smart agriculture.

REFERENCES

[1] G. Anthon, D. M. Barrett, "Standardization of a rapid spectrophotometric method for lycopene analysis," Acta Horticulture, pp. 758, 2007.

[2] D. M. Barrett, G. E. Anthon, "Color quality of tomato products," ACS Symposium Series, vol. 983, 2008, pp. 131–139.

[3] K. A. Canene, J. K. Campbell, S. Zaripheh, E. H. Jeffery, J. W. Erdman, "The tomato as a functional food," The Journal of Nutrition, vol. 135, 2005, pp. 1126-1130. doi:10.1093/jn/135.5.1226.

[4] A. Clément, R. Bacon, S. Sirois, M. Dorais, "Mature-ripe tomato spectral classification according to lycopene content and fruit type by visible, NIR reflectance and intrinsic fluorescence," Quality Assurance and Safety of Crops & Foods, vol. 7(5), 2015, pp.747-756. doi:10.3920/QAS2014.0521

[5] A. Clément, M. Dorais, M. Vernon, "Nondestructive measurement of fresh tomato lycopene content and other physicochemical characteristics using visible-NIR spectroscopy," Journal of Agricultural and Food Chemistry, vol. 56(21), 2008, pp. 9813–9818. doi:10.1021/jf801299r.

[6] A. R. Davis, W.W. Fish, P. V. Perkins, "A rapid spectrophotometric method for analyzing lycopene content in tomato and tomato products," Postharvest Biology and Technology, vol. 28, 2003, pp.425–430. doi:10.1016/S0925-5214(02)00203-X.

[7] A. R. Davis, W. W Fish, P. V. Perkins, "A rapid hexane-free method for analyzing lycopene content in watermelon," Journal of Food Science, vol. 68(1), 2003, pp. 328–332. doi:10.1111/j.1365-2621.2003.tb14160.x.

[8] A. Ford, A. Roberts, "Colour space conversions," Retrieved July 10, 2018, http://www.poynton. com/PDFs/coloureq.pdf

[9] B. Ghatak, S. B. Ali, N. Debabhuti, P. Sharma, A. Ghosh, B. Tudu, N. Bhattacharya, "Discrimination of tomatoes based on lycopene using cyclic voltammetry," Sensor Letters, vol. 15(10), 2017, pp. 827–836. doi:10.1166/ sl.2017.3884.

[10] R. M. H. Nguyen, D. K. Prasad, M. S. Brown, "Training-based spectral reconstruction from a single RGB image," Computer Vision-ECCV2014, Part VII, LNCS 8695, 2014, pp. 186–201.

[11] C. Limberg, A. Melnik, A. Harter, H. Ritter, "YOLO—You Only Look," 10647 Times. arXiv 2022, arXiv:2201.06159, Available online: https://arxiv.org/abs/2201.06159.

[12] J. S., Mommoh, J. L., Obetta, S. N., John, K. Okokpujie, O. N., Omoruyi, and A. A., Awelewa, "Detection of fruit ripeness and defectiveness using Convolutional Neural Networks," Inf. Dyn. Appl., vol. 3, no. 3, pp. 184–199, 2024. https://doi.org/10.56578/ida030304.

[13] D. Minagawa, J. Kim, "Prediction of Harvest Time of Tomato Using Mask R-CNN," AgriEngineering, vol. 4, 2022, pp. 356–366.

[14] M. Waseem, M. M. Sajjad, L H. Naqvi, Y. Majeed, T. U. Rehman, T. Nadeem, "Deep learning model for precise and rapid prediction of tomato maturity based on image recognition," Food Physics, vol. 2, 2025, 100060, https://doi.org/10.1016/j.foodp.2025.100060.

[15] Y. Okabe, T. Hiraguri, K. Endo, T. Kimura, D. Hayashi, "Classification of Tomato Harvest Timing Using an AI Camera and Analysis Based on Experimental Results," AgriEngineering, vol. 7 (2), 2025, pp. 48, https://doi.org/10.3390/agriengineering7020048.

[16] X. Ye, T. Izawa, S. Zhang, "Rapid determination of lycopene content and fruit grading in tomatoes using a smart device camera," Cogent Engineering, vol. 5(1), 2018, 15004499, DOI: 10.1080/23311916.2018.1504499.

[17] A. Kweon, V. Hu, J. Y. Lim, T. Gee, E. Liu, H. Williams, B. A. MacDonald, M. Nejati, I. Sa, H. S. Ahn, "Visual-Based Tomato Size Measurement System for an Indoor Farming," Environment, April 2023. DOI:10.48550/arXiv.2304.06177.

[18] Md. S. Alam, Md. R. Ali, A. Rahman, "Detection and localization of ripe tomato in greenhouse environment using Keras-based deep learning models," Journal of Agriculture and Food Research, vol. 23, 2025, 102182, https://doi.org/10.1016/j.jafr.2025.102182.