



Deposited via The University of Leeds.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/228653/>

Version: Accepted Version

---

**Article:**

Guo, Z., Xiao, H., Dai, Z. et al. (2025) Identification of apple variety using machine vision and deep learning with Multi-Head Attention mechanism and GLCM. *Journal of Food Measurement and Characterization*. ISSN: 2193-4126

<https://doi.org/10.1007/s11694-025-03385-5>

---

This is an author produced version of an article published in *Journal of Food Measurement and Characterization*, made available under the terms of the Creative Commons Attribution License (CC-BY), which permits unrestricted use, distribution and reproduction in any medium, provided the original work is properly cited.

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.

1           **Identification of apple variety using machine vision and deep**  
2           **learning with multi-head attention mechanism and GLCM**

3           Zhiming Guo<sup>a, b\*</sup>, Haidi Xiao<sup>a</sup>, Zhiqiang Dai<sup>c</sup>, Chen Wang<sup>d</sup>, Chanjun Sun<sup>a</sup>,  
4           Nicholas Watson<sup>b</sup>, Megan Povey<sup>b</sup>, Xiaobo Zou<sup>a</sup>

5           <sup>a</sup> *China Light Industry Key Laboratory of Food Intelligent Detection & Processing,*  
6           *School of Food and Biological Engineering, Jiangsu University, Zhenjiang 212013,*  
7           *China*

8           <sup>b</sup> *School of Food Science and Nutrition, University of Leeds, Leeds LS2 9JT, United*  
9           *Kingdom*

10          <sup>c</sup> *Beijing Key Laboratory of Multi-dimension & Multi-scale Computational*  
11          *Photography, LUSTER LightTech Co., Ltd. Beijing 100094, China*

12          <sup>d</sup> *Key Laboratory of Modern Agricultural Equipment and Technology, Ministry of*  
13          *Education, School of Agricultural Engineering, Jiangsu University, Zhenjiang 212013,*  
14          *China*

15          \*Corresponding author at: School of Food and Biological Engineering, Jiangsu  
16          University, Zhenjiang 212013, China. Email addressed: guozhiming@ujs.edu.cn (Z.  
17          Guo)

18 **Abstract:** Apple variety identification plays a crucial role in pomology and agricultural  
19 sciences, as it could effectively assist growers in optimizing orchard management,  
20 enhancing product quality, and meeting consumer demand. Traditional identification  
21 methods based on visual observation are often influenced by various factors, including  
22 human subjective judgment and inter-cultivar variability. To address these challenges,  
23 with the support of the China Agriculture Research Systems for Apple Industry and  
24 Jiangsu University, we collected sample images of eleven common apple varieties in  
25 China, followed by image enhancement and dataset expansion to establish an apple  
26 sample database. Subsequently, Convolutional Neural Network (CNN), MobileNet  
27 Version 2 (MobileNetV2), and Visual Geometry Group 19 (VGG19) neural network  
28 models were utilized for apple variety classification using image-based data.  
29 Additionally, two optimization techniques, namely Multi-Head Attention and Gray-  
30 Level Co-occurrence Matrix (GLCM), were incorporated to further improve  
31 classification accuracy. Results demonstrated that the baseline CNN achieved an  
32 accuracy of 96.46%, while MobileNetV2 and VGG19 reached 97.78% and 97.25%,  
33 respectively. Multi-Head Attention improved feature extraction but sometimes reduced  
34 performance, as observed in MobileNetV2 (87.33%). In contrast, GLCM significantly  
35 improved model accuracy, with MobileNetV2 achieving the highest accuracy (98.25%)  
36 and the lowest Mean Absolute Error (MAE) (0.0571). GLCM consistently  
37 outperformed other techniques across all models, proving particularly effective for  
38 texture-rich image classification. These findings suggest that GLCM is a powerful  
39 enhancement for deep learning models, improving accuracy, precision, and recall in  
40 apple variety classification, with MobileNetV2 combined with GLCM yielding the best  
41 overall results.

42 **Keywords:** Apple variety classification; Deep learning; Optimization techniques;  
43 Convolutional neural network

44 **List of abbreviations**

<b>Abbreviation</b>	<b>Full Term</b>
CNN	Convolutional Neural Network
MobileNetV2	Mobile Network Version 2
VGG19	Visual Geometry Group 19-layer
MAE	Mean Absolute Error
GI	Geographical Indication
CNNs	Convolutional Neural Networks
RH	Relative Humidity
TP	True Positive
TN	True Negatives
FP	False Positive
FN	False Negative
t-SNE	t-distributed stochastic neighbor embedding
PC	principal components
SSL	Self-Supervised Learning

---

## 46 **1. Introduction**

47 As one of the most widely cultivated fruits globally, apples hold substantial  
48 economic value, contributing billions of dollars annually to economies around the  
49 world. China, as one of the largest apple-producing regions, harvested 47.57 million  
50 tons of apples in 2023, accounting for more than half of the global apple production [1].  
51 The long shelf life of apples and their suitability for various preservation techniques,  
52 such as refrigeration and canning, further enhance their economic importance.

53 Apples play a pivotal role in futures markets by facilitating price discovery and  
54 risk management, allowing producers, traders, and consumers to lock in prices and  
55 mitigate risks associated with price fluctuations. Accurate identification of apple  
56 varieties is fundamental for improving price discovery in futures markets, minimizing  
57 fraud during delivery, and enhancing agricultural efficiency.

58 With the growing promotion of Geographical Indication (GI) products, the  
59 connection between specific apple varieties and their production regions has become  
60 crucial for enhancing commercial value. GI products, protected by intellectual property  
61 laws, are associated with high quality and authenticity, fostering consumer trust.  
62 Accurate identification of apple varieties is vital in preventing fraud, such as the  
63 substitution or mixing of similar varieties, which compromises the integrity of GI  
64 products. Ensuring only the correct cultivar is labeled and sold protects the authenticity  
65 of GI apples and prevents misleading claims. For example, misrepresenting a premium  
66 cultivar like "Yanfu" with a lower-quality one like "ordinary Red Fuji" damages the

67 market value and reputation of GI apples. Therefore, accurate cultivar identification is  
68 essential for improving trade transparency, promoting GI products, and protecting  
69 market integrity, thereby supporting the sustainable growth of the apple industry  
70 through fair competition and consumer trust [2].

71 However, identifying apple varieties is a significant challenge for farmers, traders,  
72 and consumers due to the large number of varieties that share similar appearances,  
73 especially within certain color ranges [3, 4]. Additionally, factors such as growing  
74 conditions, soil types, and cultivation practices influence the shape, color, and texture  
75 of apples, making the identification of apple varieties more complex [5]. In China,  
76 various apple varieties, such as Fuji, Red Delicious, and Gala, are cultivated for specific  
77 market segments. However, the introduction of new apple varieties closely resembling  
78 their parent varieties has made accurate classification challenging. Traditional methods,  
79 such as analyzing leaf and fruit characteristics, consulting experts, or using genetic  
80 testing, are time-consuming and struggle to balance efficiency with cost-effectiveness  
81 [6].

82 Since 2012, Convolutional Neural Networks (CNNs) have emerged as a leading  
83 technology in image processing and computer vision. Originally introduced by LeCun  
84 in the 1980s, CNNs serve as the foundational architecture in deep learning and have  
85 been widely applied to image classification and object detection tasks. The specific  
86 architecture is shown in Fig. 1. With local connectivity and weight sharing, CNNs  
87 enable efficient feature extraction and robust performance under varying conditions. In

88 recent years, deep learning driven by CNNs has increasingly been utilized to address  
89 complex challenges in agricultural systems [7, 8]. Applications include detecting  
90 cucumber powdery mildew [9], classifying the maturity stages of custard apple fruits  
91 through image processing [10], and predicting fruit size and weight in apples using  
92 RGB-D cameras [11]. Additionally, CNNs have been employed for target detection,  
93 with models developed for object category recognition [12, 13].

94 Apple cultivar recognition involves extracting discriminative features from visual  
95 attributes such as shape, texture, and color [14-16]. A study [17] used VGG16, VGG19,  
96 and MobileNet to distinguish between ten apple varieties, with DenseNet201 achieving  
97 97.48% accuracy. Another study [18] combined CNNs with a convolutional  
98 autoencoder to classify 26 fruits, including nine apple varieties. Additionally, a separate  
99 study [19] developed a shallow CNN to simplify deep neural networks for apple image  
100 recognition, achieving 92% accuracy.

101 For real-time applications and deployment on mobile or embedded devices,  
102 models with lower computational complexity are essential [20, 21]. While architectures  
103 like EfficientNet and Vision Transformers offer high performance, their high  
104 computational demands make them unsuitable for resource-constrained environments.  
105 Similarly, YOLO's emphasis on object localization limits its effectiveness in fine-  
106 grained classification, such as distinguishing subtle apple cultivar differences [22, 23].  
107 In contrast, CNNs, including MobileNetV2 and VGG19, provide robust feature  
108 extraction. MobileNetV2, optimized for mobile environments, reduces computational

109 costs through depth-wise separable convolutions and an inverted residual structure  
110 while maintaining strong representational power [24]. VGG19, developed by the Visual  
111 Geometry Group at Oxford in 2014, features a deeper architecture with 3×3  
112 convolutional kernels, enabling better multi-level feature extraction and enhanced  
113 capability for distinguishing subtle morphological differences in apple varieties [25].  
114 Both models balance efficiency, feature extraction, and robustness, making them ideal  
115 for apple cultivar recognition in practical, resource-limited settings [26, 27].

116 To improve model performance, two optimization techniques—Multi-Head  
117 Attention mechanism and GLCM—were used. Multi-Head Attention mechanism from  
118 the Transformer architecture enhances feature representation by capturing global  
119 dependencies [28]. GLCM, a texture analysis method, identifies pixel intensity co-  
120 occurrence patterns, aiding in the differentiation of similar apple varieties [29]. Its low  
121 computational complexity makes it ideal for lightweight models [25]. Together, these  
122 techniques improve accuracy and robustness in apple cultivar recognition.

123 The research approach is illustrated in Fig. 2. Our study utilizes a comprehensive  
124 dataset comprising eight varieties from the Fuji series, along with three additional  
125 widely cultivated apple varieties from major apple-producing regions in China. To  
126 address the challenge of new varieties closely resembling their parent varieties, we  
127 employ traditional CNN models, MobileNetV2 (efficiency), and VGG19 (feature  
128 extraction), incorporating the Multi-Head Attention mechanism and GLCM to improve  
129 feature extraction and classification accuracy. The objective is to enhance the accuracy

130 and efficacy of apple cultivar classification by developing more efficient deep learning  
131 models and optimizing CNN architectures, evaluating the interaction of various  
132 components in real-world applications, and providing a comprehensive classification  
133 model for common apple varieties in China.

## 134 **2. Material and methods**

### 135 **2.1 Apple samples**

136 According to the 2023 data from the National Bureau of Statistics of China, the  
137 major apple varieties cultivated in China include the Fuji series, Red Delicious series,  
138 Gala series, and Golden Delicious series. Within the Fuji series, sub-varieties such as  
139 Red Fuji, Yanfu, and Miyakiji are widely cultivated. Fuji apples, including these sub-  
140 varieties, account for 69.8% of the total apple cultivation area in China, highlighting  
141 their dominance in the industry [30]. For this study, eleven apple varieties were selected,  
142 including eight from the Fuji series and three additional widely cultivated varieties,  
143 representing the major apple types in China.

144 **With support from the China Agriculture Research Systems for Apple Industry,**  
145 **samples of eleven apple varieties were collected (Fig. 3). To enhance the**  
146 **generalizability of the apple sample images, samples of the same cultivar were**  
147 **sourced from different regions, improving the applicability of the apple**  
148 **classification model across diverse production areas and testing environments.**  
149 **The apple samples were gathered from major apple-producing regions across**  
150 **seven provinces in China and transported to the Apple Testing Center at Jiangsu**  
151 **University. Upon arrival, the apples were stored in the cold storage of the**  
152 **laboratory at a temperature of 0-2°C and a relative humidity (RH) of 85-90%.2.2**

### 153 **Image acquisition**

154       The image collection for all evaluated apple varieties was completed within three  
155 days of receiving the samples at the Apple Testing Center at Jiangsu University. These  
156 samples were sourced from cooperatives, companies, and experimental stations in  
157 major apple-producing regions across China. The image collection process took place  
158 from September 18, 2023, to January 28, 2024.

159       Upon receipt of the samples, we classified them according to the GB/T 10651-  
160 2008 Chinese National Standard for Fresh Apples, categorizing them into Extra Class,  
161 Class I, Class II, and Substandard [31]. To minimize the impact of surface defects on  
162 feature extraction, apples classified as Class II or above were selected for image  
163 acquisition. The selection criteria were as follows: 1) A fruit diameter of at least 65 mm;  
164 2) A shape index of 0.8 or higher; 3) A surface coloration rate of over 30%; 4) No more  
165 than two surface defects; 5) A total damage area smaller than one square centimeter.

166       Advanced hardware was utilized in the image acquisition process to ensure high-  
167 quality data collection. A BVC8350LC, a 3CMOS color area scan camera (Blue Vision

168 Corporation, Japan) served as the primary imaging device. This industrial-grade camera  
169 featured a resolution of 324 million pixels and was equipped with an f/1.2 maximum  
170 aperture lens to capture fine details. Full-spectrum industrial lighting was employed to  
171 provide uniform and consistent illumination, minimizing the influence of external  
172 lighting variations. The camera was equipped with a fixed-focus lens with a focal length  
173 of 20 mm and operated at a working distance of 650 mm. It operated at a shutter speed  
174 of 1/60 s to ensure sharp image capture. All images were saved in BMP format to  
175 preserve their quality and facilitate subsequent processing.

176 An average of 50 apple samples from each cultivar were selected for image  
177 collection. During this process, each apple sample was placed on a white-background  
178 platform for image capture, with images taken at 90-degree intervals along the apple's  
179 equator. Additionally, two images focusing on the apple's stem and calyx were captured,  
180 resulting in a total of six images per apple. Detailed information about the eleven  
181 photographed apple varieties was shown in Fig. 3. The dataset was then divided into  
182 training, validation, and test sets in a 6:2:2 ratio for model training.

### 183 **2.3 Database construction**

184 The original apple images were preprocessed to enhance data diversity and expand  
185 the dataset, thereby mitigating overfitting during model training. To reduce the  
186 influence of irrelevant information and highlight the apple as the primary subject, the  
187 images, initially with a resolution of 2448×1840 pixels, were cropped to 1324×1256  
188 pixels. Four augmentation techniques—darkening, variation, Gaussian filtering, and

189 Gaussian noise—were applied to the images. For each cultivar, 50 images were  
190 randomly selected from the set and processed with each technique, ensuring a total of  
191 approximately 500 sample images per cultivar, as shown in Table 1.

## 192 **2.4 Experimental design and algorithms**

### 193 *2.4.1. Experimental design*

194 **Image preprocessing:** The quality of the apple images was affected by  
195 environmental factors, such as low light and vibration, as well as varietal characteristics.  
196 Preprocessing steps included cropping to focus on the apples, color correction for  
197 consistent lighting, noise reduction, brightness and contrast adjustment, and  
198 normalization. These measures improved image quality, ensuring reliable model  
199 training.

200 **Model Training:** Three base models (CNN, MobileNetV2, and VGG19) were  
201 employed for training. Additionally, GLCM features and Multi-Head Attention  
202 mechanisms were integrated with these models to create composite architectures. This  
203 combination harnesses both low-level texture information and high-level semantic  
204 features, enhancing classification accuracy and robustness. In total, nine distinct deep  
205 learning models were developed during the training process.

206 **Model Evaluation and Deployment:** Accuracy, precision, recall, and MAE were  
207 commonly used parameters for evaluating model performance [32]. The trained models  
208 were tested for prediction to identify issues related to overfitting or underfitting [2].  
209 The validated models were uploaded to the server, and the local device accessed the  
210 server through software developed using PyCharm Professional (version 2023.1,

211 JetBrains, Prague, Czech Republic) and Qt Creator (version 5.0.2, the Qt Company,  
212 Espoo, Finland) on a system running Windows 10 (version 22H2, Microsoft, Redmond,  
213 WA, USA) with Python (version 3.7.6, Python Software Foundation, Wilmington, DE,  
214 USA). This software was designed for downloading the models to perform image  
215 acquisition, preprocessing, and cultivar prediction.

#### 216 ***2.4.2. Implementation description***

217 The research was conducted on a Windows 10 (version 22H2, Microsoft,  
218 Redmond, WA, USA) operating system with Python (version 3.7.6, Python Software  
219 Foundation, Wilmington, DE, USA) and the TensorFlow framework, leveraging a Tesla  
220 P4 GPU (NVIDIA Corporation, Santa Clara, CA, USA) for computation. These  
221 computational resources ensured efficient task execution within the projected  
222 timeframe.

223 The apple cultivar recognition system was developed using TensorFlow (version  
224 2.3.0, TensorFlow, Inc., Mountain View, CA, USA), a widely adopted Python deep  
225 learning framework. To maintain experimental rigor and fairness, a concise ten-layer  
226 CNN architecture, comprising convolutional and pooling layers, was employed. This  
227 design balanced simplicity and performance, minimizing overfitting risks. Both the  
228 MobileNetV2 and VGG19 architectures were integrated into the framework,  
229 maintaining uniform training parameters. Each model underwent preliminary training  
230 for 10 epochs to assess stability and determine the optimal number of training steps.  
231 The final training was conducted in 70 epochs, which accurately reflected the model's

232 validation accuracy and loss throughout the process. The Multi-Head Attention  
233 mechanism employed consists of eight attention heads, with a hidden feature dimension  
234 of 256. The Adam optimizer was adopted for all models, with cross-entropy loss serving  
235 as the loss function to optimize classification performance.

### 236 ***2.4.3. Evaluation metrics***

237 The confusion matrix is an essential tool for evaluating classification model  
238 performance and for assessing multi-class models. It provides a clear representation of  
239 predicted versus actual outcomes for each class, highlighting classification  
240 misclassifications [33]. The confusion matrix primarily consists of parameters such as  
241 accuracy, precision, recall, and F1 score, facilitating a comprehensive assessment of the  
242 model's performance.

243 These metrics depend on four fundamental values derived from the confusion  
244 matrix: true positives (TP), true negatives (TN), false positives (FP), and false negatives  
245 (FN). Specifically, precision represents the ratio of true positives to all predicted  
246 positives, while recall indicates the proportion of true positives among all actual  
247 positives. The F1 score, calculated as the harmonic mean of precision and recall,  
248 effectively balances these two metrics. It is particularly advantageous in scenarios with  
249 class imbalance, as it considers both the precision and recall of correct positive  
250 predictions.

251 In addition, MAE is a regression metric that quantifies the average absolute  
252 difference between predicted and actual values. It is computed as the mean of the

253 absolute residuals over all samples, offering an intuitive and interpretable measure of  
 254 prediction accuracy [34].

255 These performance indicators can be mathematically defined as follows:

256

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

$$Precision = \frac{TP}{TP + FP} \quad (2)$$

$$Recall = \frac{TP}{TP + FN} \quad (3)$$

$$F1\ Score = \frac{2 \times TP}{2 \times TP + FP + FN} \quad (4)$$

$$MAE = \frac{1}{n} \quad (5)$$

257 Moreover, five-fold cross-validation was employed to assess the model's  
 258 performance, ensuring robustness given the sample size and specific experimental  
 259 requirements. This approach divided the dataset into five subsets, with each subset  
 260 serving as a validation set once, while the others were utilized for training. By  
 261 systematically rotating through these subsets, this approach mitigated the risks of  
 262 overfitting and underfitting, thereby improving the model's predictive performance.  
 263 Furthermore, this robust evaluation method provided a comprehensive assessment of  
 264 the model's generalization ability on unseen data, ensuring stable and reliable  
 265 performance in real-world applications.

## 266 **3. Results and discussion**

### 267 **3.1. Model training**

#### 268 *3.1.1 Model performance evaluation*

269 The training dataset comprised approximately 5,500 images, encompassing a  
270 diverse range of apple varieties widely cultivated in China, including the Fuji series and  
271 newer varieties such as Venus Gold and Yuhua Fushi. Because of the large size of the  
272 raw apple images, all images were resized to 224×224 pixels. The t-SNE (t-distributed  
273 stochastic neighbor embedding) dimensionality reduction of the training samples was  
274 shown in Fig. 4. Before performing t-SNE, PCA was applied to reduce the  
275 dimensionality to 50 components, and the first two principal components (PC1 and PC2)  
276 explained 41.36% of the variance (PC1: 27.10%, PC2: 14.26%). Most apple varieties  
277 exhibited significant confusion in feature distribution, such as Yanfu No.10 and Yanfu  
278 No.3. Conversely, varieties such as Holstein and Venus Gold exhibited clearer  
279 clustering due to distinct differences in their color features compared to other varieties.  
280 In contrast, Nagafu No.2, was more dispersed across the clusters of other varieties,  
281 suggesting that it posed a greater challenge for identification.

282 To evaluate the performance of the CNN models, training accuracy and loss rate  
283 graphs were generated and illustrated in Fig. 5.

284 The CNN series models exhibited relatively smooth fluctuations overall in Fig.  
285 5(a). However, significant variability was observed in the later stages of the CNN model  
286 incorporating Multi-Head Attention, which suggested overfitting due to excessive  
287 exposure to certain features. This was confirmed by its higher overall loss rate. Similar

288 trends were observed in the MobileNetV2 and VGG19 models, highlighting the impact  
289 of the Multi-Head Attention mechanism on performance.

290 All three foundational models were found to perform well when integrating  
291 features with GLCM in Fig. 5(a) and Fig. 5(b). Among them, the combination of  
292 MobileNetV2 and GLCM demonstrated the best overall performance, as evidenced by  
293 a stable accuracy curve and a closely following loss rate curve. Additionally, integrating  
294 GLCM with the VGG19 model significantly reduced the variability of the validation  
295 loss curve, indicating that GLCM effectively mitigated the overfitting issues observed  
296 with this model.

### 297 ***3.1.2. Comparative analysis of the overall performance of deep learning models***

298 The performance of CNN, MobileNetV2, and VGG19 model architectures for  
299 apple cultivar classification was evaluated, with the test results presented in Table 2.  
300 Classification accuracy, precision, recall, F1 score, and MAE were selected as the  
301 evaluation metrics for these models.

302 The baseline CNN model achieved excellent performance with an accuracy of  
303 96.46%, precision of 96.48%, recall of 96.46%, and F-score of 96.41%. The MAE for  
304 this model was 0.1311, which is relatively high compared to other optimized models  
305 [35, 36]. Upon incorporating the Multi-Head Attention mechanism into the CNN model,  
306 the accuracy slightly increased to 97.08%, accompanied by improvements in precision,  
307 recall, and F1-score. The precision reached 97.10%, the recall was 97.08%, and the F-  
308 score was 97.04%. The MAE decreased to 0.1090, reflecting a reduction in error. This

309 suggests that the inclusion of Multi-Head Attention enabled the model to focus on more  
310 relevant features, thereby enhancing its classification performance.

311 On the other hand, integrating GLCM with the CNN model resulted in even better  
312 performance. The accuracy increased to 97.92%, with precision, recall, and F-score all  
313 showing similar improvements, reaching 97.91%, 97.92%, and 97.86%, respectively.  
314 The MAE further decreased to 0.0980, demonstrating that the integration of GLCM not  
315 only enhanced the model's classification accuracy but also significantly reduced errors.  
316 The GLCM method, which extracts texture features from images, likely enabled the  
317 CNN to better capture subtle visual patterns specific to the apple varieties.

318 The MobileNetV2 model, known for its computational efficiency, also  
319 demonstrated strong performance. In its baseline configuration, the model achieved an  
320 accuracy of 97.78%, with both precision and recall at 97.78%, and an F-score of 97.74%.  
321 The MAE was 0.0696, indicating the best error performance among the baseline models  
322 tested. When Multi-Head Attention was applied to MobileNetV2, the accuracy  
323 decreased to 87.33%, with precision, recall, and F-score following a similar decline.  
324 The MAE increased dramatically to 0.5490. This suggests that, in this case, the Multi-  
325 Head Attention mechanism did not provide the expected improvements and negatively  
326 impacted the model's performance. Such complexity may have interfered with the  
327 added complexity from the attention mechanism interfered with MobileNetV2's  
328 efficient feature extraction, resulting in overfitting or poor generalization.

329 Integrating GLCM with MobileNetV2 resulted in a significant improvement. The

330 accuracy increased to 98.25%, with precision and recall both reaching 98.29% and  
331 98.25%, respectively. The F-score was 98.20%, while the MAE dropped to 0.0571.  
332 GLCM's texture features enhanced MobileNetV2's performance, leading to  
333 improvements in both classification accuracy and error reduction.

334 The VGG19 model, a deeper network known for its efficient feature extraction  
335 capabilities, achieved an accuracy of 97.25%, with both precision and recall at 97.25%,  
336 and an F-score of 97.22%. The MAE was 0.0992, indicating acceptable performance.  
337 However, when Multi-Head Attention was incorporated, the accuracy slightly  
338 decreased to 97.01%, with corresponding declines in precision, recall, and F1-score.  
339 The MAE decreased to 0.0975, suggesting that while the attention mechanism  
340 contributed to error reduction, it did not substantially enhance classification  
341 performance.

342 Consistent with the findings from other models, integrating GLCM with VGG19  
343 led to improved performance. The accuracy reached 97.92%, with both precision and  
344 recall at 97.91%, and the F-score at 97.66%. The MAE decreased to 0.0921. These  
345 results demonstrate that the texture features extracted by GLCM enhanced VGG19's  
346 classification performance and reduced prediction errors.

347 To conclude, the application of GLCM demonstrated consistent improvements in  
348 classification performance across all models, resulting in significant advancements in  
349 accuracy and error reduction. The incorporation of Multi-Head Attention, however,  
350 produced mixed results, improving some models, such as CNN, but significantly

351 degrading MobileNetV2's performance. Overall, GLCM was a valuable optimization  
352 technique, enhancing classification accuracy and reducing MAE, particularly for  
353 MobileNetV2 and CNN.

### 354 **3.2. Model comparison of algorithms within the same series**

355 Heatmap visualizations confirmed that models integrating GLCM and the Multi-  
356 Head Attention mechanism, built upon CNN, MobileNetV2, and VGG19 architectures,  
357 effectively focused on relevant image regions during feature extraction. These  
358 visualizations demonstrated that each model achieved satisfactory classification  
359 performance across most apple varieties.

#### 360 ***3.2.1. CNN series cultivar detection models***

361 The prediction results for eleven varieties of apples were shown in Fig. 6 from  
362 three models: (a) CNN model, (b) CNN+Multi-Head Attention model, and (c)  
363 CNN+GLCM model. The CNN model demonstrated generally high accuracy, with  
364 categories such as Changhong (95.93%) and Red Fuji (95.07%) achieving high  
365 classification accuracy. Other categories, such as Chengji No.1, Holstein, Lifu No.2,  
366 Miyakuj, Venus Gold, and Yuhua Fushi, also achieved near 100% accuracy, indicating  
367 solid performance in distinguishing these apple varieties. However, there were  
368 exhibited considerable misclassification, such as Nagafu No.2 (76.52%) and Yanfu  
369 No.3, which showed lower accuracy and higher misclassification rates. This suggested  
370 that the CNN model struggled to distinguish certain varieties which had visually similar  
371 features.

372 The addition of Multi-Head Attention improved the model's performance,  
373 particularly in categories such as Changhong (97.29%) and Red Fuji (96.06%), where  
374 accuracy improved compared to the CNN model. It also enhanced the overall prediction  
375 reliability for categories such as Chengji No.1 and Holstein, which now exhibited 100%  
376 accuracy. The Multi-Head Attention mechanism appears to have enabled the model to  
377 focus more effectively on key features, enhancing its classification performance,  
378 particularly in complex scenarios. However, despite the improved performance over the  
379 CNN model, some categories, such as Nagafu No.2 (79.13%), still showed confusion,  
380 though the performance was better than that of the CNN model.

381 The CNN + GLCM model demonstrated a significant performance improvement,  
382 particularly for categories such as Changhong (97.74%) and Red Fuji (95.07%),  
383 achieving high accuracy consistently. The integration of GLCM enhanced the model's  
384 ability to capture texture features, enabling it to more effectively differentiate between  
385 varieties with similar visual characteristics. The confusion for Nagafu No.2 (76.52%)  
386 persisted, but the accuracy improved slightly compared to the CNN model, indicating  
387 that GLCM contributed to better differentiation of such varieties. Overall, the  
388 combination of CNN and GLCM resulted in improved performance, particularly for  
389 texture-based features, although some categories still exhibited minor  
390 misclassifications.

391 In conclusion, the CNN model demonstrated strong classification performance,  
392 which was further enhanced by integrating Multi-Head Attention and GLCM, resulting

393 in improved accuracy, especially for categories with similar visual features. However,  
394 some categories, such as Nagafu No.2, still posed classification challenges.

### 395 ***3.2.2. Analysis of MobileNetV2 series cultivar detection models***

396 The original MobileNetV2 model demonstrated strong performance in predicting  
397 varieties, particularly for Changhong, Chengji No.1, and Holstein, achieving accuracies  
398 of 99.55%, 99.46%, and 100%, respectively (Fig. 7). However, the model struggled  
399 with varieties such as Nagafu No.2, Red Fuji, and Yanfu No.10, with accuracies ranging  
400 from 0% to 5%. The model excelled in identifying distinct varieties but struggled with  
401 those with less distinctive features, particularly Nagafu No.2. Despite these challenges,  
402 the model effectively distinguished the most distinctive apple varieties but faced  
403 difficulty classifying more complex or visually similar ones.

404 Integrating Multi-Head Attention into the MobileNetV2 model, its performance  
405 improved, particularly on more challenging varieties. This enhancement helped the  
406 model better capture intricate patterns, particularly for varieties such as Nagafu No.2  
407 and Red Fuji, where accuracy increased. While the model's performance on Changhong  
408 decreased slightly to 97.74%, varieties such as Yanfu No.10 and Yanfu No.3 showed  
409 notable improvements, with accuracies of 91.79% and 83.33%, respectively.  
410 Nevertheless, the model continued to face challenges with certain varieties, such as  
411 Miyakuj and Venus Gold, where misclassifications persisted.

412 The final version, which combined MobileNetV2 with GLCM, showed the most  
413 significant improvements in classification accuracy. This hybrid model performed  
414 exceptionally well, achieving 100% accuracy on Changhong, 99.46% on Chengji No.1,

415 and 98.98% on Holstein. It also excelled in identifying varieties such as Red Fuji  
 416 (97.44%) and Venus Gold (100%). “The GLCM-based approach, focusing on texture  
 417 features, helped differentiate visually similar varieties more effectively, improving  
 418 precision, especially for challenging varieties like Yanfu No.3 (96.15%).

419       The combination of MobileNetV2 with the Multi-Head Attention mechanism and  
 420 GLCM resulted in noticeable performance variations. While MobileNetV2 with GLCM  
 421 demonstrated strong classification abilities, the inclusion of Multi-Head Attention  
 422 resulted in a decline in overall performance. This decline was especially noticeable in  
 423 the identification of varieties such as Nagafu No.2, Red Fuji, and Yanfu No.3.  
 424 MobileNetV2, designed for efficiency with depthwise separable convolutions, is  
 425 lightweight, but the addition of Multi-Head Attention, which captures global context,  
 426 and increased computational complexity. This added complexity likely hindered  
 427 performance, particularly on smaller datasets, and may have contributed to overfitting.

428       The inclusion of GLCM notably enhanced classification performance, particularly  
 429 for varieties such as Nagafu No.2 and Yanfu No.3. By extracting textural features,  
 430 GLCM improved MobileNetV2's ability to capture subtle texture variations that might  
 431 have been overlooked by the model, leading to better classification accuracy.

### 432       ***3.2.3. Analysis of VGG19 series cultivar detection models***

433       VGG19 exhibited robust performance in classifying apple varieties, consistently  
 434 achieving high accuracy, as illustrated in Fig. 8. The model did not exhibit overfitting  
 435 for most varieties, achieving prediction accuracies of 98.64% for Changhong and 99.46%  
 436 for Chengji No. 1, demonstrating excellent performance. Holstein reached 98.98%

437 accuracy, while Miyakuj attained 99.48%, further showcasing the robustness of VGG19.  
438 Although the classification accuracy for Lifu No. 2 slightly decreased to 96.02%, it  
439 remained high overall. However, Nagafu No. 2 showed a relatively lower accuracy of  
440 81.74%, revealing challenges in recognizing this cultivar. Despite this, VGG19  
441 continued to perform strongly across most varieties, handling classification tasks with  
442 only minor difficulties for a few specific cases.

443       The addition of the Multi-Head Attention mechanism significantly enhanced the  
444 performance of VGG19 in classifying apple varieties. While the accuracy for  
445 Changhong decreased slightly to 97.32%, it remained high. Chengji No. 1 and Holstein  
446 maintained perfect classification at 100%, and Miyakuj also achieved 100%,  
447 highlighting the effectiveness of the Multi-Head Attention mechanism. Lifu No. 2  
448 showed stability with an accuracy of 99.50%. However, Nagafu No. 2 experienced a 4%  
449 drop in accuracy to 74.00%, indicating that the mechanism might have introduced  
450 interference for some varieties. Other varieties, such as Red Fuji and Venus Gold, saw  
451 modest improvements, with accuracies of 98.70% and 98.00%, respectively. The  
452 classification accuracy for Yanfu No. 10 and Yanfu No. 3 remained mostly unchanged  
453 at 98.50% and 92.50%. Yuhua Fushi continued to perform flawlessly with an accuracy  
454 of 100%.

455       After incorporating GLCM into the VGG19 model, the performance for many  
456 varieties showed improvement. Nagafu No.2 saw a significant increase to 84.17%, up  
457 from 81.74% in the original VGG19 model. The model's ability to correctly identify

458 "Red Fuji" remained consistent at 98.22%. Yanfu No.3 experienced a modest  
459 improvement to 94.06%, with only 0.60% misclassified. Meanwhile, varieties such as  
460 Chengji No.1 (100%) and Holstein (100%) showed no major changes, as the addition  
461 of texture-based features from GLCM did not impact these already high-performing  
462 categories. This indicates that GLCM had a substantial effect on the more challenging  
463 varieties.

464       The performance of the VGG19 model changed notably with the addition of Multi-  
465 Head Attention and GLCM, highlighting the interplay between these mechanisms and  
466 the network's structure. VGG19's deep architecture, with its multiple convolutional  
467 layers and pooling operations, effectively captured hierarchical image features.  
468 However, the introduction of Multi-Head Attention sometimes caused interference,  
469 particularly with varieties such as Nagafu No. 2, where the model struggled to  
470 distinguish subtle differences. This may have been due to an overemphasis on less  
471 relevant features. On the other hand, GLCM improved texture-based feature extraction,  
472 boosting Nagafu No. 2's accuracy. This improvement likely stemmed from GLCM's  
473 ability to capture subtle texture differences that VGG19 alone might have missed,  
474 which was especially beneficial for varieties with more intricate surface textures where  
475 color alone was insufficient for accurate classification.

476       GLCM provides complementary information about the spatial relationship  
477 between pixels, helping the model discern finer details that distinguish varieties such  
478 as Nagafu No. 2. This texture-based enhancement aids the model in focusing on

479 important patterns that might be overlooked in color-based feature extraction, resulting  
480 in improved classification accuracy for varieties with complex textures.

### 481 **3.3 Model comparison of algorithms across different series**

482       Among the nine models assessed, MobileNetV2+GLCM achieved the highest  
483 performance, achieving an overall classification accuracy of 98.25% with an MAE of  
484 0.0571. It was followed by VGG19+GLCM, which achieved a classification accuracy  
485 of 97.92% and an MAE of 0.0921. The training-related charts show that the  
486 introduction of GLCM positively impacted all three baseline models, significantly  
487 enhancing their performance. For instance, the inclusion of GLCM improved the  
488 MobileNetV2 model's accuracy to 98.25%, outperforming the baseline model.  
489 Similarly, the CNN model saw an increase in accuracy to 97.92%, maintaining strong  
490 performance compared to the version without optimization techniques.

491       In contrast, the Multi-Head Attention mechanism exhibited a selective effect. After  
492 its introduction, slight improvements were observed in the performance of CNN, with  
493 accuracy rising from 96.46% to 97.08%. However, the Multi-Head Attention  
494 mechanism had a pronounced negative impact on the MobileNetV2 model, leading to  
495 a decline in performance. A slight decrease in performance was also noted when Multi-  
496 Head Attention was applied to the VGG19 model. The confusion matrix for  
497 MobileNetV2 with Multi-Head Attention revealed a significant drop in prediction  
498 accuracy across several varieties, resulting in distorted predictions. Notable  
499 misclassifications occurred, particularly among closely related varieties such as Yanfu

500 No. 10 and Yanfu No. 3. This suggested that Multi-Head Attention did not effectively  
501 enhance MobileNetV2's performance, potentially causing overfitting or feature loss.

502 To analyze why the Multi-Head Attention mechanism had adverse effects on  
503 MobileNetV2, we considered both the model's parameter count and the nature of the  
504 attention mechanism. MobileNetV2 is a lightweight model with only 2,263,108  
505 parameters, far fewer than CNN (53,767,748) and VGG19 (26,585,163). Designed for  
506 efficiency, MobileNetV2 was optimized to work with limited computational resources,  
507 making it highly effective for tasks that require fewer parameters. However, the  
508 introduction of Multi-Head Attention, which adds complexity and increases  
509 computational demands, may have disrupted this balance, negatively impacting  
510 performance.

511 The introduction of the Multi-Head Attention mechanism added computational  
512 complexity and increased the model's capacity, which may have led to overfitting or  
513 instability, especially in lightweight models such as MobileNetV2. These effects were  
514 particularly noticeable in cases where the dataset or training process was sensitive to  
515 parameter adjustments. The observed training fluctuations and the decline in  
516 recognition accuracy indicated that the added complexity disrupted MobileNetV2's  
517 balance, ultimately reducing its efficiency and performance. In contrast, models with  
518 higher parameter capacities, such as CNN, better accommodated the additional layers  
519 introduced by attention mechanisms, which explains their improved performance.

520 Finally, we observed that the recognition accuracy for the Nagafu No. 2 cultivar

521 remained around 80% across the nine models, with frequent misclassifications as Lifu  
522 No. 2, Miyakuj, or Venus Gold. Similarly, the recognition accuracy for Yanfu No. 3  
523 clustered around 92%, indicating a high degree of similarity between Yanfu No. 3 and  
524 other apple varieties. The persistent misclassification of Nagafu No. 2, a key maternal  
525 parent in Fuji-lineage breeding, is likely attributable to its similar genomic traits, which  
526 lead to phenotypic ambiguities in standard RGB imaging and, consequently, contribute  
527 to classification errors.

528       To further improve cultivar recognition accuracy, plan to implement data  
529 augmentation techniques such as rotation, scaling, flipping, and color adjustments.  
530 These methods will introduce variability, increase dataset diversity, and enhance the  
531 model's generalization capability. Additionally, the loss function will be modified by  
532 incorporating class weights to place greater emphasis on reducing misclassifications of  
533 underrepresented varieties.

534       We will also integrate additional features, such as color histograms, shape  
535 descriptors, and metadata, to better differentiate between varieties. Furthermore, to  
536 enhance the model's feature extraction capabilities, we will explore advanced  
537 techniques such as Adaptive Feature Fusion and Self-Supervised Learning (SSL).  
538 These methods will enable more effective high-level feature learning, ultimately  
539 improving the model's accuracy in recognizing apple varieties with genetic  
540 relationships, such as Nagafu No. 2.

## 541 **4. Conclusion**

542 This research presented an innovative approach to classifying and recognizing  
543 apple varieties using deep learning techniques, investigating the integration of the  
544 Multi-Head Attention mechanism and GLCM optimization across three distinct  
545 architectures: CNN, MobileNetV2, and VGG19. By combining image processing with  
546 machine learning, the study significantly improved the accuracy and efficiency of apple  
547 cultivar identification. Focusing on eleven popular apple varieties in China, the  
548 optimized models consistently achieved classification accuracies above 95%, with  
549 some varieties exceeding 98%.

550 MobileNetV2+GLCM achieved the highest accuracy of 98.25%, demonstrating  
551 the effectiveness of combining traditional methods with advanced image processing  
552 techniques. Compared to the study in [17], which trained seven types of CNN models  
553 on a dataset comprising 5,808 images from 10 different Turkish apple varieties and  
554 identified DenseNet as the best-performing model with an accuracy of 97.48%, and the  
555 study in [37], which employed MobileNetV2 and EfficientNetV2B0 to classify six  
556 Turkish apple varieties from a dataset of 120 images, where EfficientNetV2B0  
557 combined with GLCM and Color-Space achieved the highest accuracy of 98.33%, our  
558 research covered a broader range of apple varieties and utilized a significantly larger  
559 dataset. Moreover, the MobileNetV2 model offers advantages over EfficientNet and  
560 DenseNet in terms of computational efficiency and suitability for deployment in  
561 resource-constrained environments.

562           However, integrating Multi-Head Attention into lightweight models such as  
563 MobileNetV2 resulted in training instability and potential overfitting. In contrast, its  
564 incorporation into CNN and VGG19 resulted in moderate improvements, with  
565 accuracies exceeding 97.08%, demonstrating the selective advantages of advanced  
566 feature enhancement[38, 39]. These findings provide valuable insights for developing  
567 models for fruit species recognition and surface defect detection in agricultural products  
568 such as citrus, pears, and peaches. By enhancing the precision and efficiency of  
569 agricultural technology, these models have the potential to be applied across a wide  
570 range of agricultural applications, from quality control to automated sorting.

571           Our research focuses on identifying 11 apple varieties commonly cultivated in  
572 China, contributing to the development of deep learning models for apple recognition.  
573 The study holds significant potential for integration into commercial apple sorting lines.  
574 These models can complement existing systems for quality recognition, thereby  
575 enhancing both sorting efficiency and accuracy. Furthermore, the system is d designed  
576 for deployment on mobile devices, enabling farmers, traders, and consumers to  
577 conveniently use the software and models for real-time cultivar identification and  
578 quality assessment.

579           Despite these advantages, the apple images used in our study were collected under  
580 controlled conditions using a professional acquisition platform. In practical applications,  
581 factors such as lighting, vibrations, and the presence of foreign objects may affect image  
582 quality. To improve model robustness and generalizability, it is essential to incorporate

583 images captured in real-world conditions, such as those obtained from conveyor belts,  
584 storage environments, and varying lighting scenarios. Additionally, apple phenotypes  
585 vary with maturity, causing certain varieties to exhibit striking similarities at specific  
586 growth stages while diverging significantly at others, which poses challenges for  
587 classification, sorting, and grading models.

588 To further validate the efficacy of the proposed approach, future research should  
589 expand both the dataset and the diversity of varieties to enhance broader applicability.  
590 We will focus on optimizing the baseline models to improve robustness, building upon  
591 the current findings. Additionally, we will explore the incorporation of additional  
592 features, such as color, shape, and temporal information, to enhance grading accuracy.  
593 Furthermore, real-time applications will be developed to improve automation and  
594 efficiency in apple cultivar recognition systems.

## 595 **Acknowledgments**

596 This research was funded by the Key R&D Project of Jiangsu Province  
597 (BE2022363), the National Key R&D Program of China (2022YFD2100604), the  
598 National Natural Science Foundation of China (W2412103), the Natural Science  
599 Foundation of Jiangsu Province (BK20220524), and the Senior Talent Program of  
600 Jiangsu University (22JDG015).

## References

- Fan, S., Liang, X., Huang, W., Zhang, V. J., Pang, Q., He, X., Li, L., & Zhang, C. 2022. Real-time defects detection for apple sorting using NIR cameras with pruning-based YOLOV4 network. *Comput. Electron. Agric.*, 193, Article 106715. <https://doi.org/10.1016/j.compag.2022.106715>
- Musacchi, S., & Serra, S. 2018. Apple fruit quality: Overview on pre-harvest factors. *Scientia Hort.*, 234, 409-430. <https://doi.org/10.1016/j.scienta.2017.12.057>
- Wu, D., Lv, S., Jiang, M., & Song, H. 2020. Using channel pruning-based YOLO v4 deep learning algorithm for the real-time and accurate detection of apple flowers in natural environments. *Comput. Electron. Agric.*, 178, Article 105742. <https://doi.org/10.1016/j.compag.2020.105742>
- Adhinata, F. D., Wahyono, & Sumiharto, R. 2024. A comprehensive survey on weed and crop classification using machine learning and deep learning. *Artif. Intell. Agric*, 13, 45-63. <https://doi.org/10.1016/j.aiia.2024.06.005>
- Liang, X., Zhang, R., Gleason, M. L., & Sun, G. 2022. Sustainable Apple Disease Management in China: Challenges and Future Directions for a Transforming Industry. *Plant Dis.*, 106(3), 786-799. <https://doi.org/10.1094/pdis-06-21-1190-fe>
- Gao, F., Fu, L., Zhang, X., Majeed, Y., Li, R., Karkee, M., & Zhang, Q. 2020. Multi-class fruit-on-plant detection for apple in SNAP system using Faster R-CNN. *Comput. Electron. Agric.*, 176, Article 105634. <https://doi.org/10.1016/j.compag.2020.105634>
- Knott, M., Perez-Cruz, F., & Defraeye, T. 2023. Facilitated machine learning for image-based fruit quality assessment. *J. Food Eng.*, 345, Article 111401. <https://doi.org/10.1016/j.jfoodeng.2022.111401>
- Zheng, S., Gao, P., Zhang, J., Ma, Z., & Chen, S. 2024. A precise grape yield prediction method based on a modified DCNN model. *Comput. Electron. Agric.*, 225, Article 109338. <https://doi.org/10.1016/j.compag.2024.109338>
- Yag, I., & Altan, A. 2022. Artificial Intelligence-Based Robust Hybrid Algorithm Design and Implementation for Real-Time Detection of Plant Diseases in Agricultural Environments. *Biology*, 11(12), Article 1732. <https://doi.org/10.3390/biology11121732>

- Wakchaure, G. C., Nikam, S. B., Barge, K. R., Kumar, S., Meena, K. K., Nagalkar, V. J., Choudhari, J. D., Kad, V. P., & Reddy, K. S. 2024. Maturity stages detection prototype device for classifying custard apple (*Annona squamosa* L) fruit using image processing approach. *Smart Agric. Technol.*, 7, Article 100394. <https://doi.org/10.1016/j.atech.2023.100394>
- Miranda, J. C., Arno, J., Gene-Mola, J., Lordan, J., Asin, L., & Gregorio, E. 2023. Assessing automatic data processing algorithms for RGB-D cameras to predict fruit size and weight in apples. *Comput. Electron. Agric.*, 214, Article 108302. <https://doi.org/10.1016/j.compag.2023.108302>
- Wang, C., Liu, S., Wang, Y., Xiong, J., Zhang, Z., Zhao, B., Luo, L., Lin, G., & He, P. 2022. Application of Convolutional Neural Network-Based Detection Methods in Fresh Fruit Production: A Comprehensive Review. *Front. Plant Sci.*, 13, Article 868745. <https://doi.org/10.3389/fpls.2022.868745>
- Lei, L., Yang, Q., Yang, L., Shen, T., Wang, R., & Fu, C. 2024. Deep learning implementation of image segmentation in agricultural applications: a comprehensive review. *Artif. Intell. Rev.*, 57(6), Article 149. <https://doi.org/10.1007/s10462-024-10775-6>
- Ji, W., Wang, J., Xu, B., & Zhang, T. 2023. Apple Grading Based on Multi-Dimensional View Processing and Deep Learning. *Foods*, 12(11), Article 2117. <https://doi.org/10.3390/foods12112117>
- Wang, Z., Jin, L., Wang, S., & Xu, H. 2022. Apple stem/calyx real-time recognition using YOLO-v5 algorithm for fruit automatic loading system. *Postharvest Biol. Technol.*, 185, Article 111808. <https://doi.org/10.1016/j.postharvbio.2021.111808>
- Ji, W., Zhang, T., Xu, B., & He, G. 2024. Apple recognition and picking sequence planning for harvesting robot in a complex environment. *J. Agric. Eng.*, 55(1), Article 1549. <https://doi.org/10.4081/jae.2024.1549>
- Taner, A., Mengstu, M. T., Selvi, K. C., Duran, H., Gur, I., & Ungureanu, N. 2024. Apple Varieties Classification Using Deep Features and Machine Learning. *Agriculture-Basel*, 14(2), Article 252. <https://doi.org/10.3390/agriculture14020252>
- Xue, G., Liu, S., & Ma, Y. 2020. A hybrid deep learning-based fruit classification using attention model and convolution autoencoder. *Complex & Intelligent Systems*, 9, 2209–2219. <https://doi.org/10.1007/s40747-020-00192-x>

- Li, J., Xie, S., Chen, Z., Liu, H., Kang, J., Fan, Z., & Li, W. 2020. A Shallow Convolutional Neural Network for Apple Classification. *IEEE Access*, 8. <https://doi.org/10.1109/ACCESS.2020.3002882> (IEEE)
- Guo, Z., Zou, Y., Sun, C., Jayan, H., Jiang, S., El-Seedi, H. R., & Zou, X. 2024. Nondestructive determination of edible quality and watercore degree of apples by portable Vis/NIR transmittance system combined with CARS-CNN. *J. Food Meas. Charact.*, 4058–4073, (2024). <https://doi.org/10.1007/s11694-024-02476-z>
- Xu, B., Cui, X., Ji, W., Yuan, H., & Wang, J. 2023. Apple Grading Method Design and Implementation for Automatic Grader Based on Improved YOLOv5. *Agriculture-Basel*, 13(1), Article 124. <https://doi.org/10.3390/agriculture13010124>
- Wang, A., Qian, W., Li, A., Xu, Y., Hu, J., Xie, Y., & Zhang, L. 2024. NVW-YOLOv8s: An improved YOLOv8s network for real-time detection and segmentation of tomato fruits at different ripeness stages. *Comput. Electron. Agric.*, 219, Article 108833. <https://doi.org/10.1016/j.compag.2024.108833>
- Lu, S., Chen, W., Zhang, X., & Karkee, M. 2022. Canopy-attention-YOLOv4-based immature/mature apple fruit detection on dense-foliage tree architectures for early crop load estimation. *Comput. Electron. Agric.*, 193, Article 106696. <https://doi.org/10.1016/j.compag.2022.106696>
- Liu, Y., Wang, Z., Wang, R., Chen, J., & Gao, H. 2023. Flooding-based MobileNet to identify cucumber diseases from leaf images in natural scenes. *Comput. Electron. Agric.*, 213, Article 108166. <https://doi.org/10.1016/j.compag.2023.108166>
- Bansal, M., Kumar, M., Sachdeva, M., & Mittal, A. 2021. Transfer learning for image classification using VGG19: Caltech-101 image data set. *J. Ambient Intell. Humaniz. Comput.*, 14, 3609–3620. <https://doi.org/10.1007/s12652-021-03488-z>
- Sun, J., Zhang, L., Zhou, X., Yao, K., Tian, Y., & Nirere, A. 2021. A method of information fusion for identification of rice seed varieties based on hyperspectral imaging technology. *J. Food Process Eng.*, 44(9), Article e13797. <https://doi.org/10.1111/jfpe.13797>
- Thakur, P. S., Sheorey, T., & Ojha, A. 2023. VGG-ICNN: A Lightweight CNN model for crop disease identification. *Multimed. Tools Appl.*, 82(1), 497-520. <https://doi.org/10.1007/s11042-022-13144-z>

- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., & Polosukhin, I. 2017. Attention is all you need. *NeurIPS*, 30, Article 1706.03762. <https://doi.org/10.48550/arXiv.1706.03762>
- Iqbal, N., Mumtaz, R., Shafi, U., & Zaidi, S. M. H. 2021. Gray level co-occurrence matrix (GLCM) texture based crop classification using low altitude remote sensing platforms. *PeerJ Comput. Sci.*, Article e536. <https://doi.org/10.7717/peerj-cs.536>
- BEEDATA. (2024). *Yunguo: 2024 China Apple Industry Data Analysis Report*. <https://www.weihengag.com/home/article/detail/id/23912>
- (SAC), S. A. o. C. (2008). *Fresh apples (GB/T 10651-2008)*.
- Jia, W., Zhang, Z., Shao, W., Hou, S., Ji, Z., Liu, G., & Yin, X. 2021. FoveaMask: A fast and accurate deep learning model for green fruit instance segmentation. *Comput. Electron. Agric.*, 191, Article 106488. <https://doi.org/10.1016/j.compag.2021.106488>
- Jia, W., Tian, Y., Luo, R., Zhang, Z., Lian, J., & Zheng, Y. 2020. Detection and segmentation of overlapped fruits based on optimized mask R-CNN application in apple harvesting robot. *Comput. Electron. Agric.*, 172, Article 105380. <https://doi.org/10.1016/j.compag.2020.105380>
- Shafik, W., Tufail, A., Liyanage, C. D. S., & Apong, R. A. A. H. M. 2023. Using a novel convolutional neural network for plant pests detection and disease classification. *J. Sci. Food Agric.*, 103(12), 5849-5861. <https://doi.org/10.1002/jsfa.12700>
- Rai, N., Zhang, Y., Ram, B. G., Schumacher, L., Yellavajjala, R. K., Bajwa, S., & Sun, X. 2023. Applications of deep learning in precision weed management: A review. *Comput. Electron. Agric.*, 206, Article 107698. <https://doi.org/10.1016/j.compag.2023.107698>
- Pacal, I., Kunduracioglu, I., Alma, M. H., Deveci, M., Kadry, S., Nedoma, J., Slany, V., & Martinek, R. 2024. A systematic review of deep learning techniques for plant diseases. *Artif. Intell. Rev.*, 57(11), Article 304. <https://doi.org/10.1007/s10462-024-10944-7>
- Kilicarslan, S., Donmez, E., & Kilicarslan, S. 2024. Identification of apple varieties using hybrid transfer learning and multi-level feature extraction. *Eur. Food Res. Technol.*, 250(3), 895-909. <https://doi.org/10.1007/s00217-023-04436-1>
- Li, Y., Yao, T., Pan, Y., & Mei, T. 2023. Contextual Transformer Networks for Visual Recognition. *IEEE Trans. Geosci. Remote Sens.*, 45(2), 1489-1500.

<https://doi.org/10.1109/tpami.2022.3164083>

Zhou, Y., Lao, C., Yang, Y., Zhang, Z., Chen, H., Chen, Y., Chen, J., Ning, J., & Yang, N. 2021.

Diagnosis of winter-wheat water stress based on UAV-borne multispectral image texture and vegetation indices. *Agric. Water Manag.*, 256, Article 107076.

<https://doi.org/10.1016/j.agwat.2021.107076>

## Figure Captions

**Fig. 1.** Schematic of the steps involved in intelligent identification of apples based on deep learning.

**Fig. 2.** Images of the eleven major apple cultivars in the detection dataset created by Jiangsu University and supported by China Agriculture Research Systems for the Apple Industry.

**Fig. 3.** Apple classification algorithm based on deep learning, including (a) Apple dataset, (b) Multi-Head Attention mechanism, (c) Detection result, (d) CNN algorithm, (e) VGG19 algorithm, and (f) MobileNet-V2 algorithm.

**Fig. 4.** t-SNE Dimensionality Reduction Distribution of Apple Samples.

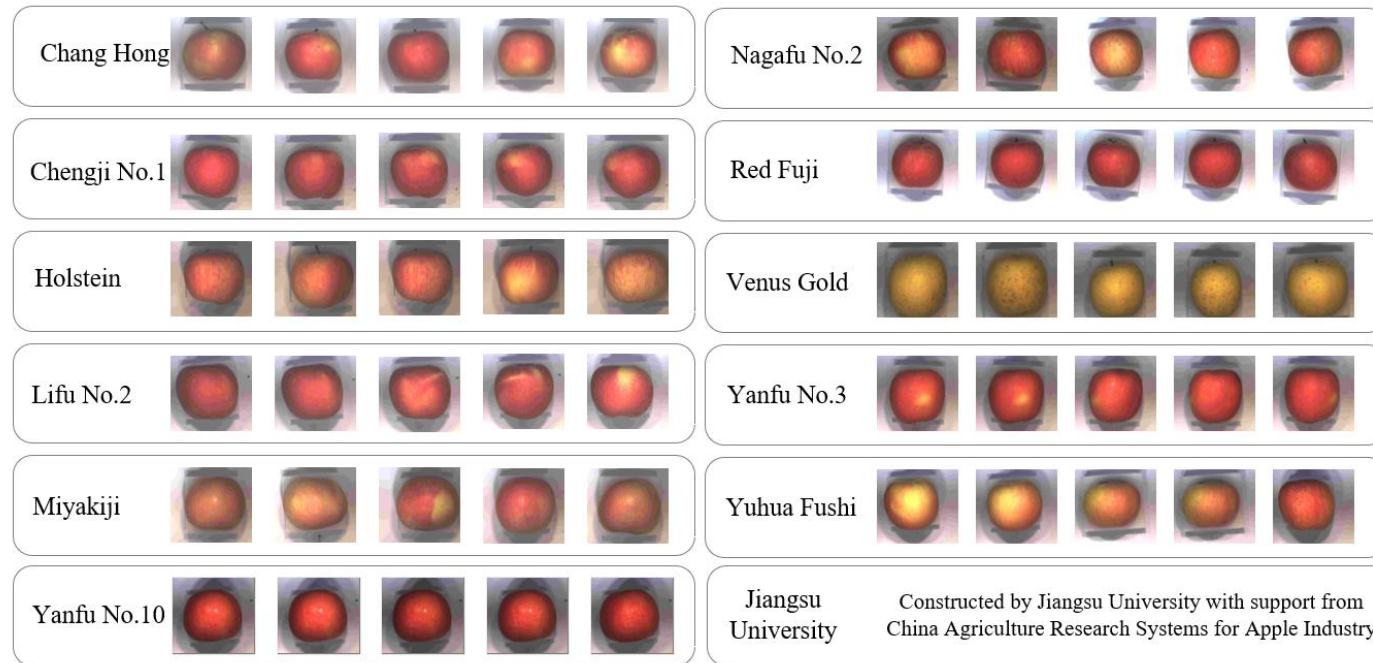
**Fig. 5.** Comparison of Training Accuracy and Validation Accuracy of the Eight Proposed Models.

**Fig. 6.** Prediction of eleven apple cultivars using three CNN-based models, including (a) CNN model, (b) CNN+Multi-Head Attention model, and (c) CNN+GLCM model.

**Fig. 7.** Prediction of eleven apple cultivars using three MobileNet V2-based models, including (a) MobileNet V2 model, (b) MobileNet V2+Multi-Head Attention model, and (c) MobileNet V2+GLCM model.

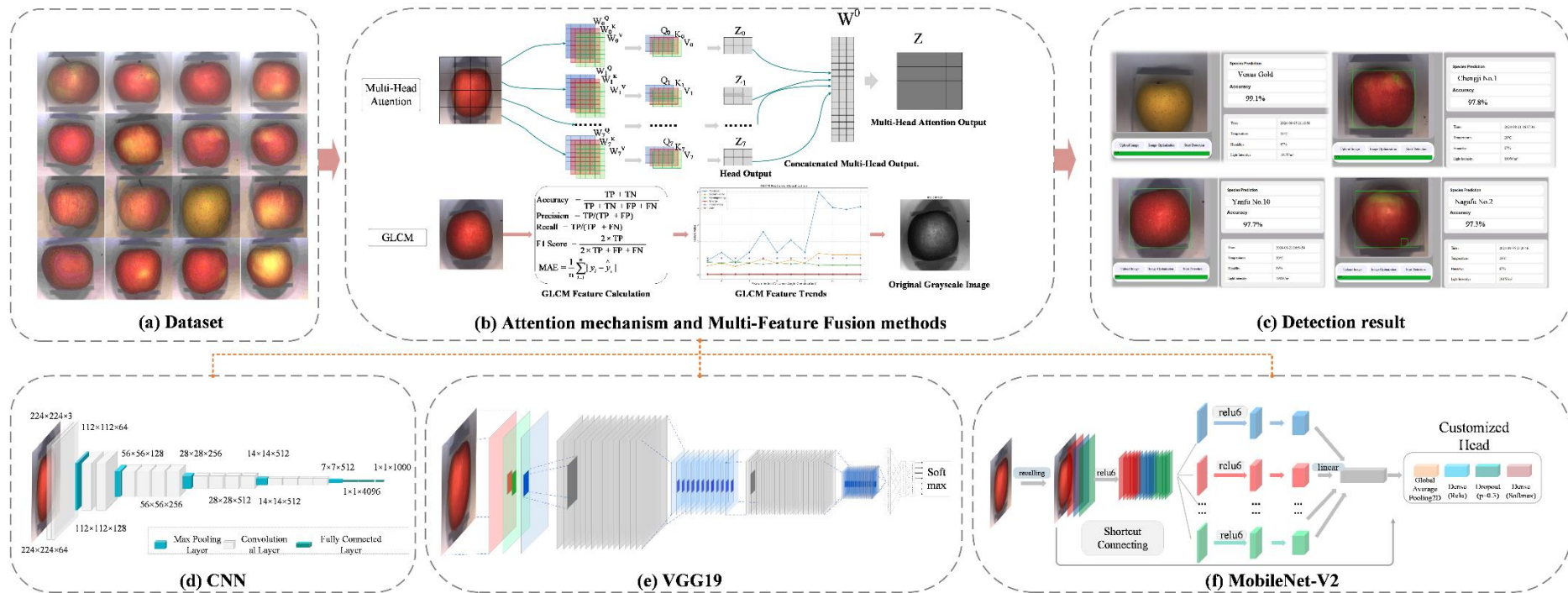
**Fig. 8.** Prediction of eleven apple cultivars using two VGG19-based models, including (a) VGG19 model, (b) VGG19+GLCM model, and (c) VGG19+Multi-Head Attention model.





**Fig. 2.** Images of the eleven major apple cultivars in the detection dataset created by Jiangsu University and supported by China Agriculture Research Systems for the Apple Industry.

## Identification of apple by machine vision and deep learning



**Fig. 3.** Apple classification algorithm based on deep learning, including (a) Apple dataset, (b) Multi-Head Attention mechanism, (c) Detection result, (d) CNN algorithm, (e) VGG19 algorithm, and (f) MobileNet-V2 algorithm.

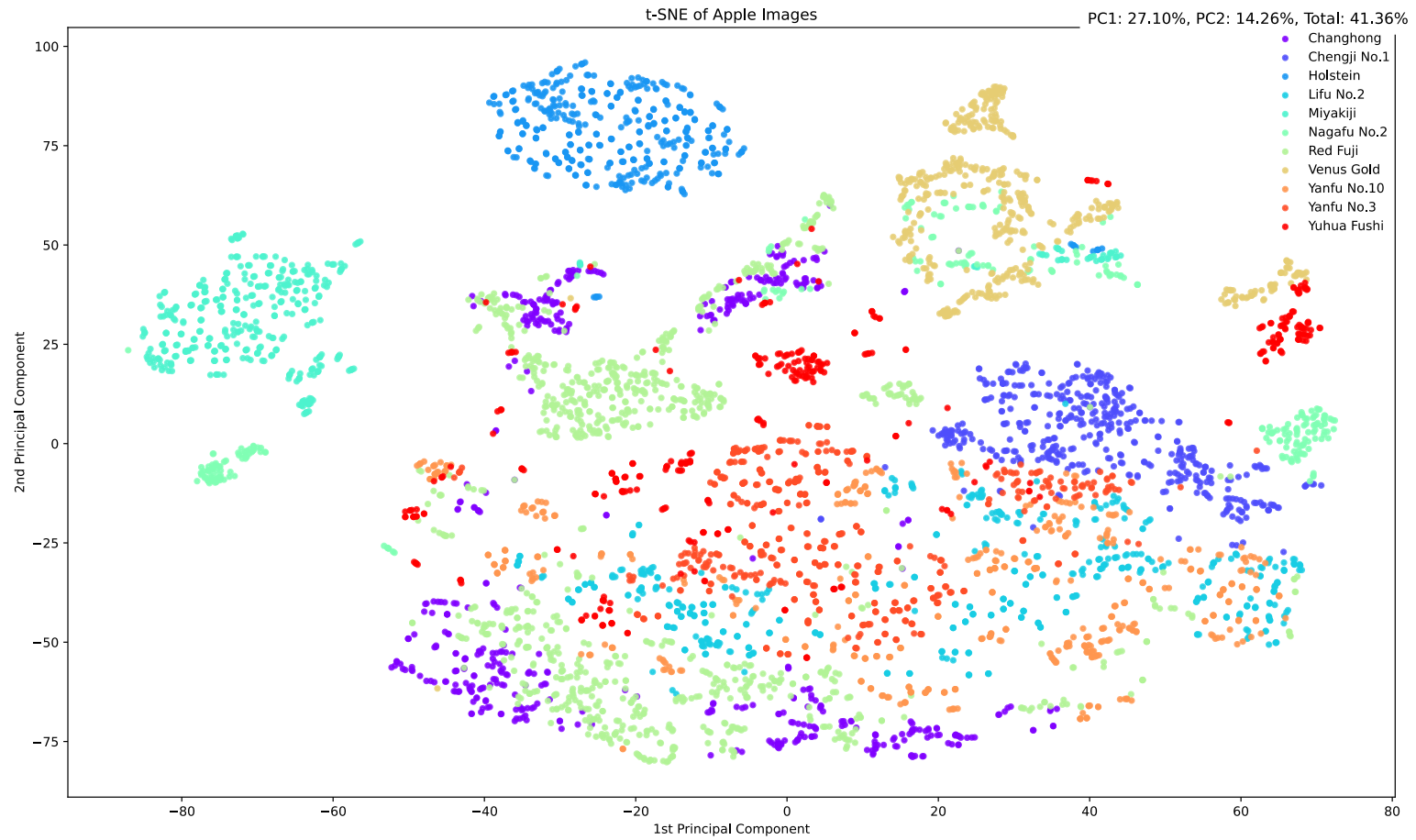
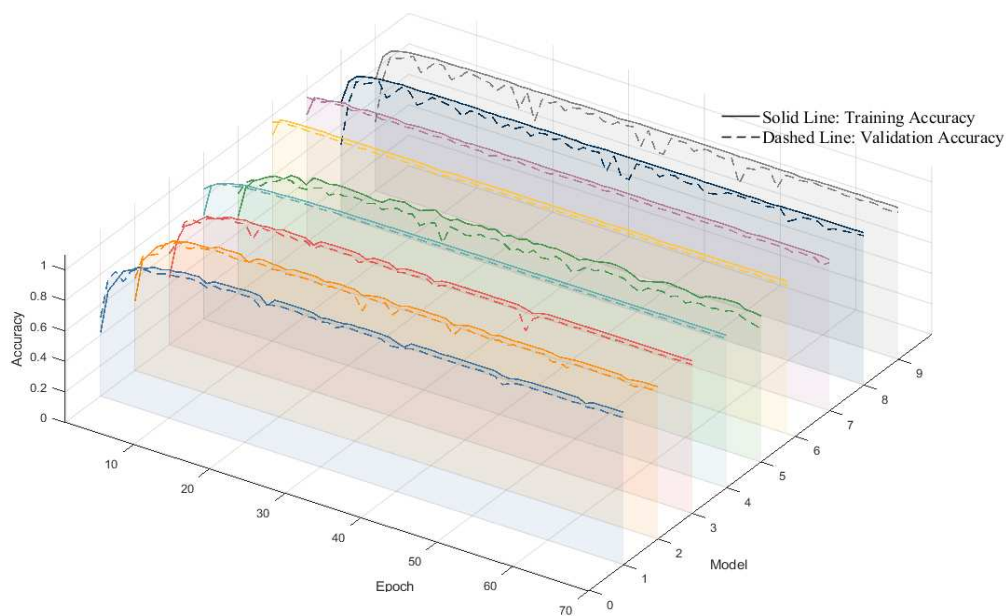


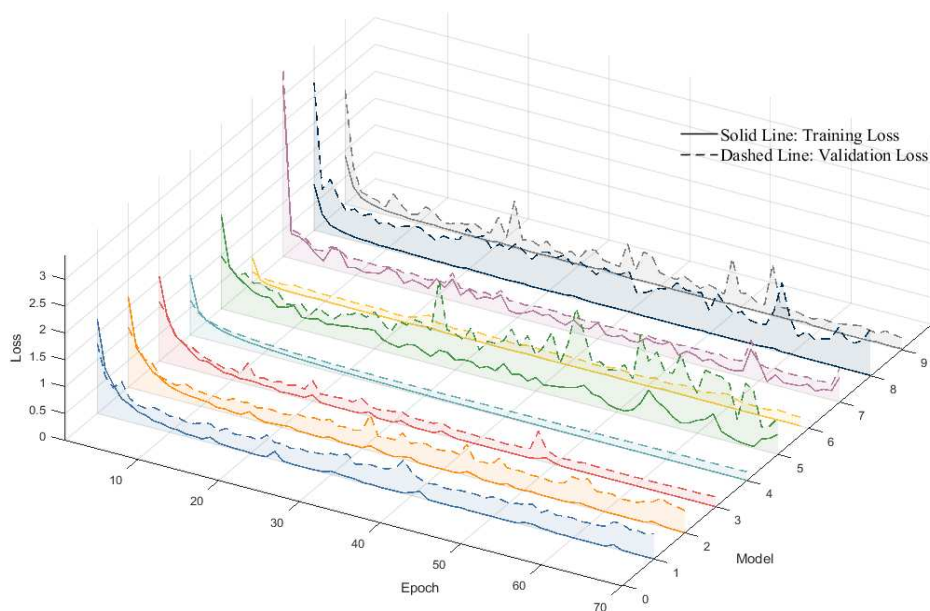
Fig. 4. t-SNE dimensionality reduction distribution of apple samples.

(a) Training Accuracy and Validation Accuracy



- Model 1: CNN
- Model 2: CNN+Multi-Head Attention
- Model 3: CNN+GLCM
- Model 4: MobileNetV2
- Model 5: MobileNetV2+Multi-Head Attention
- Model 6: MobileNetV2+GLCM
- Model 7: VGG19
- Model 8: VGG19+MultiHead Attention
- Model 9: VGG19+GLCM

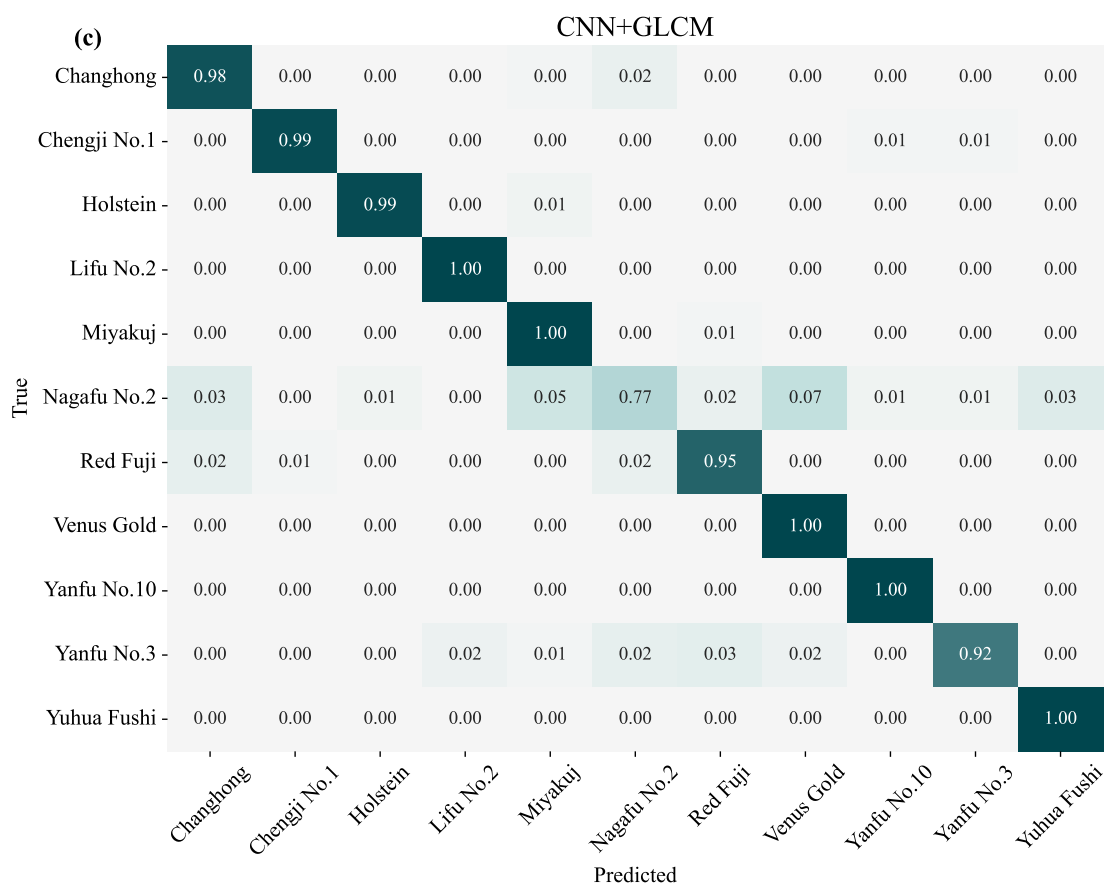
(b) Training Loss and Validation Loss



**Fig. 5.** Comparison of Training Accuracy and Validation Accuracy of the eight proposed models.

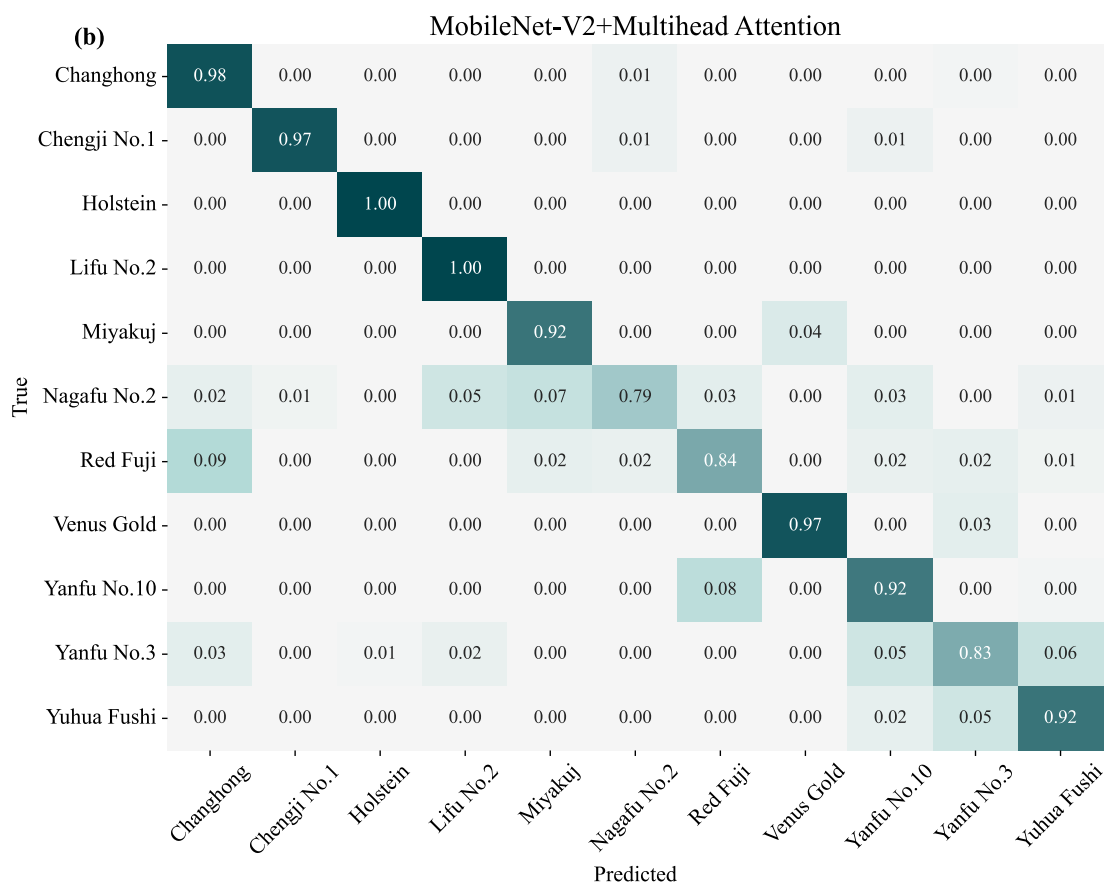
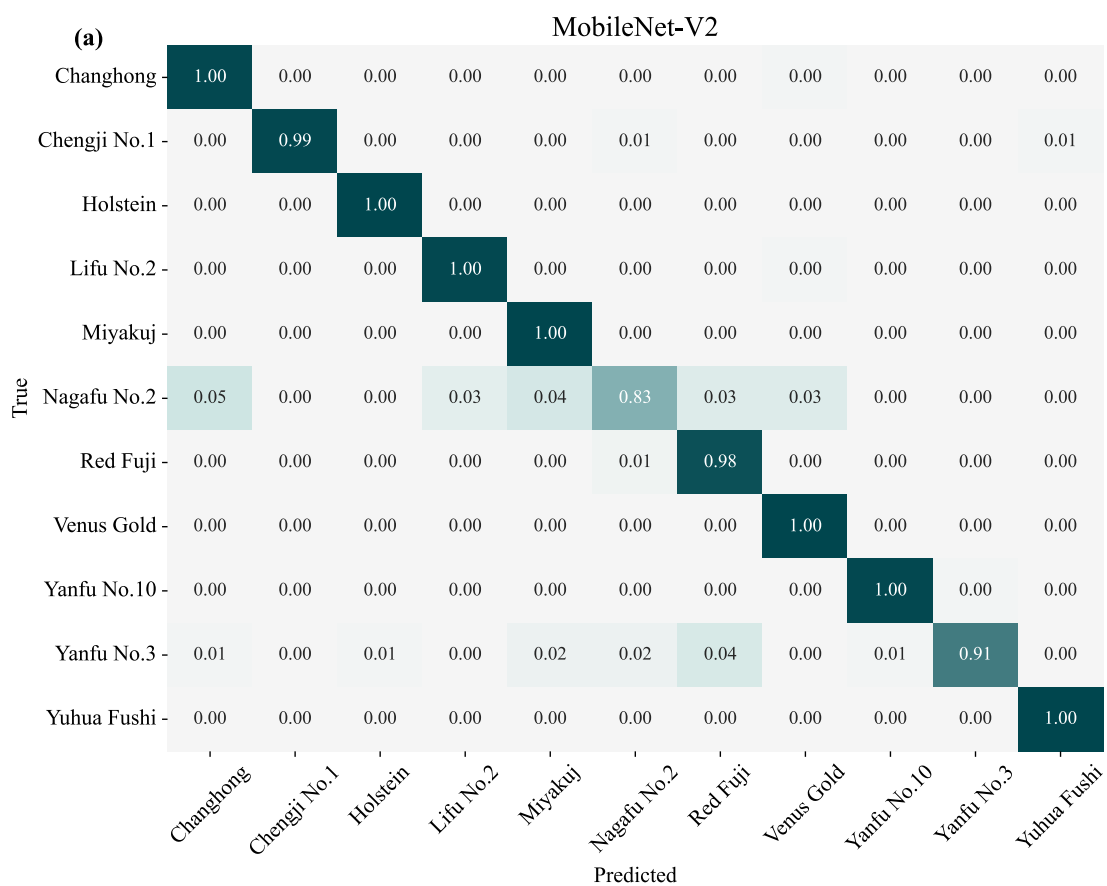


Identification of apple by machine vision and deep learning



**Fig. 6.** Prediction of eleven apple cultivars using three CNN-based models, including (a) CNN model, (b) CNN+Multi-Head Attention model, and (c) CNN+GLCM mode.

## Identification of apple by machine vision and deep learning



Identification of apple by machine vision and deep learning

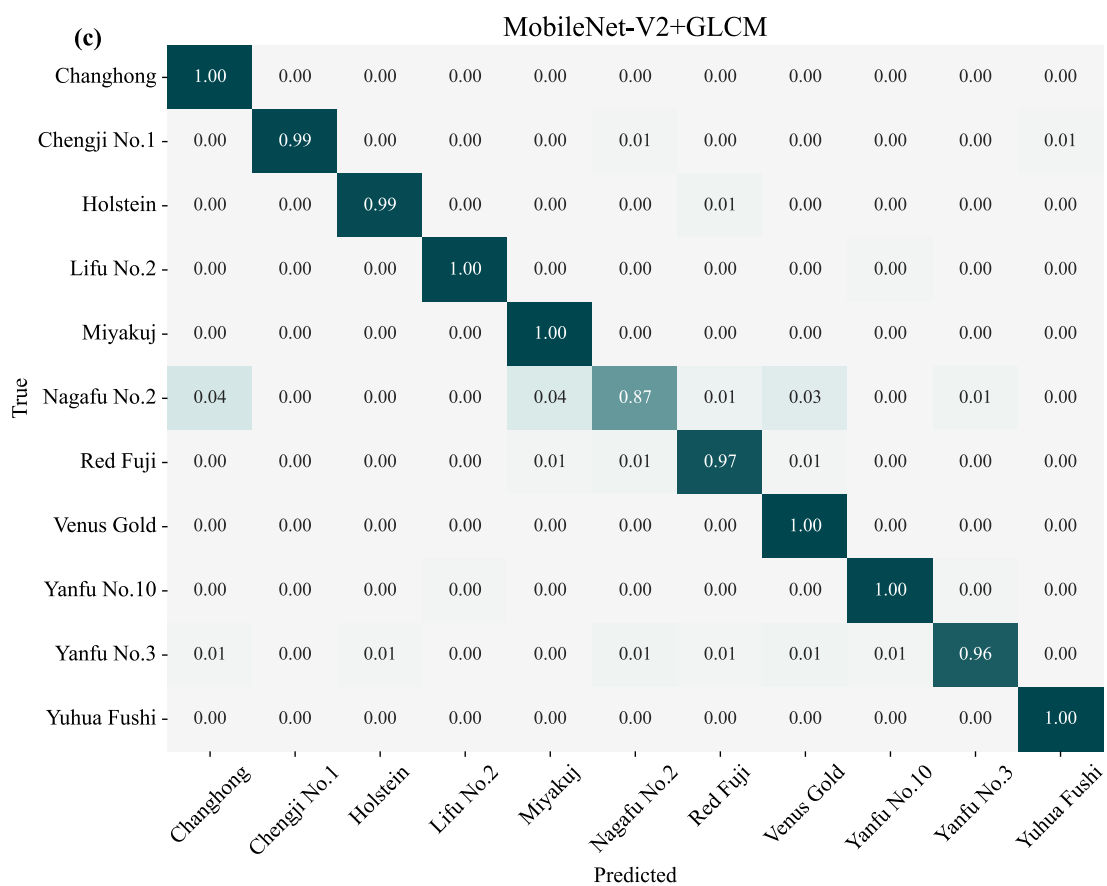
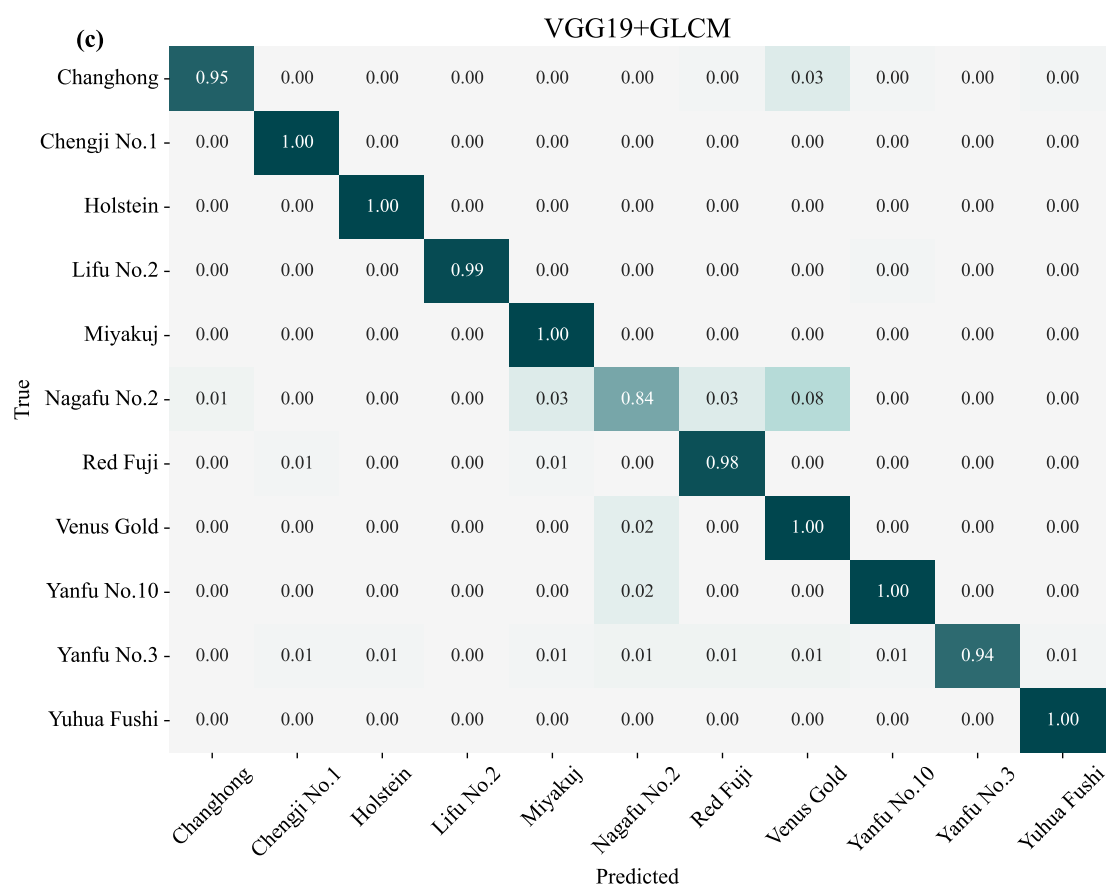


Fig. 7. Prediction of eleven apple cultivars using three MobileNet-V2-based models, including (a) MobileNet-V2 model, (b) MobileNet-V2+Multi-Head Attention model, and (c) MobileNet-V2+GLCM model.





**Fig. 8.** Prediction of eleven apple cultivars using three VGG19-based models, including (a) VGG19 model, (b) VGG19+Multi-Head Attention model and (c) VGG19+GLCM model.

## **Table Captions**

**Table 1** Apple cultivar recognition dataset, featuring initial images of eleven apple samples and their augmented versions.

**Table 2** Performance of deep learning models combining Multi-Head Attention mechanism and Gray-Level Co-occurrence Matrix with CNN, MobileNet-V2, and VGG19.

**Table 1** Apple cultivar recognition dataset, featuring original images of eleven apple cultivar samples and their augmented versions.

Cultivar Name	Apple	Original	Darkened	Varied	Gaussian	Gaussian	Total
	Origin	Images	Images	Images	Filtered	Noise	
		Images	Images	Images	Images	Images	Images
Changhong	Shaanxi Province	294	50	50	50	50	494
Chengji No.1	Shanxi Province	294	50	50	50	50	494
Holstein	Shanxi Province	312	50	50	50	50	512
Lifu No.2	Henan Province	306	50	50	50	50	506
Miyakji	Hebei Province	312	50	50	50	50	512
Nagafu No.2	Hebei Province	312	50	50	50	50	512
Red fuji	Hebei Province	288	50	50	50	50	488
Venus Gold	Liaoning Province	288	50	50	50	50	488
Yanfu No.3	Shandong Province	312	50	50	50	50	512
Yanfu No.10	Shandong Province	324	50	50	50	50	524
Yuhua Fushi	Shandong Province	300	50	50	50	50	500

**Table 2** Model performance of CNN, MobileNet-V2, and VGG19.

Morphological method	Model	Accuracy	Precision	Recall	F-score	MAE
	optimization content					
CNN	/	0.9646	0.9648	0.9646	0.9641	0.1311
CNN	Multi-Head Attention	0.9708	0.971	0.9708	0.9704	0.1090
CNN	GLCM	0.9792	0.9791	0.9792	0.9786	0.0980
MobileNet-V2	/	0.9778	0.9778	0.9778	0.9774	0.0696
MobileNet-V2	Multi-Head Attention	0.8733	0.8834	0.8733	0.8704	0.5490
MobileNet-V2	GLCM	0.9825	0.9829	0.9825	0.9820	0.0571
VGG19	/	0.9725	0.9725	0.9725	0.9722	0.0992
VGG19	Multi-Head Attention	0.9701	0.9703	0.9700	0.9698	0.0975
VGG19	GLCM	0.9792	0.9791	0.9772	0.9766	0.0921